

Applied Multivariate Analysis

Myung Hee Lee

Fall 2014

Tentative Schedule

- introduction
- descriptive techniques
- matrix algebra
- multivariate distribution
- multivariate normal distribution
- principal component analysis
- discriminant analysis
- factor analysis
- cluster analysis
- canonical correlation analysis
- graphical models

Multivariate data set?

- A data set in which several variables are measured on each sampled unit.
- We have n units (individuals) and p variables in a multivariate data set, where $p > 1$.
- Notation:

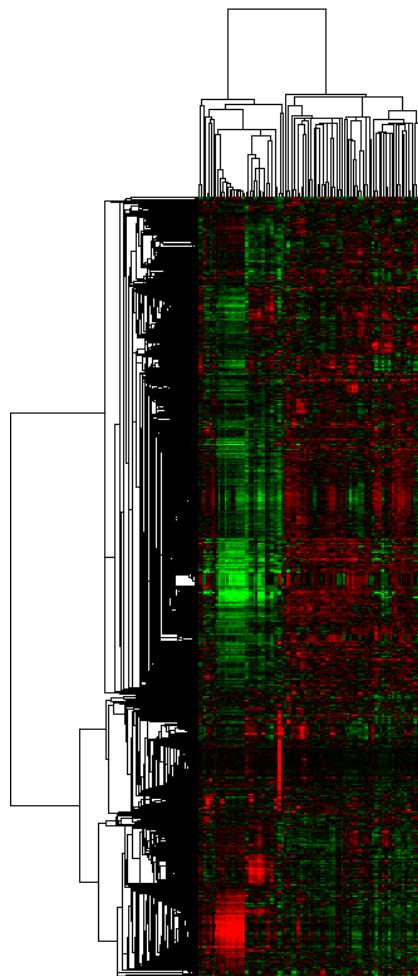
Visualization of multivariate data?

- loosely, one can think of each variables as an additional dimension
- many variable data = data sitting in a high dimensional space
- how do we look at the data?
 - 1-d? histogram
 - 2-d? scatter plot
 - 3-d? 3-D scatter plot? (spinning) point cloud view
 - 100-d?

Visualization of multivariate data?

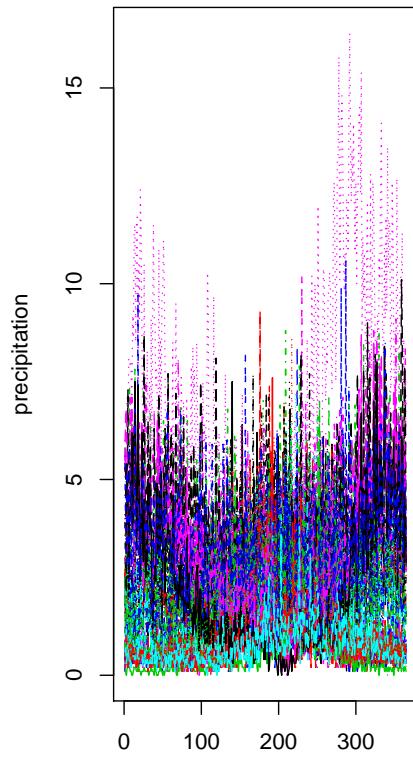
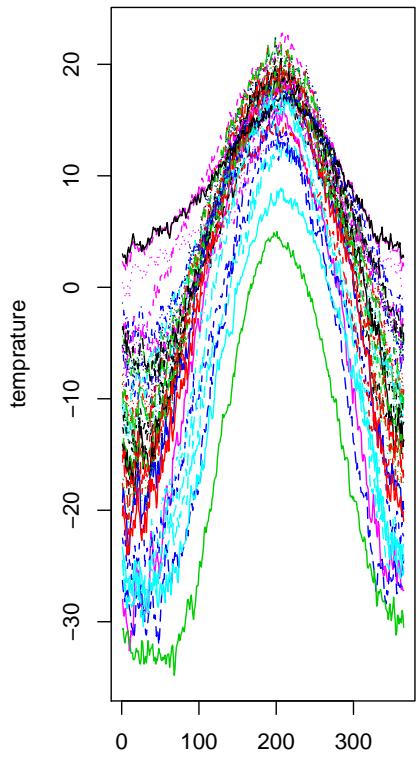
- loosely, one can think of each variables as an additional dimension
- many variable data = data sitting in a high dimensional space
- how do we look at the data?
 - 1-d? histogram
 - 2-d? scatter plot
 - 3-d? 3-D scatter plot? (spinning) point cloud view
 - 100-d?
 - projection: there are many directions of interest on which projection is useful

Example: Ovarian Cancer Data



$n = 108$ and $p = 28,475$

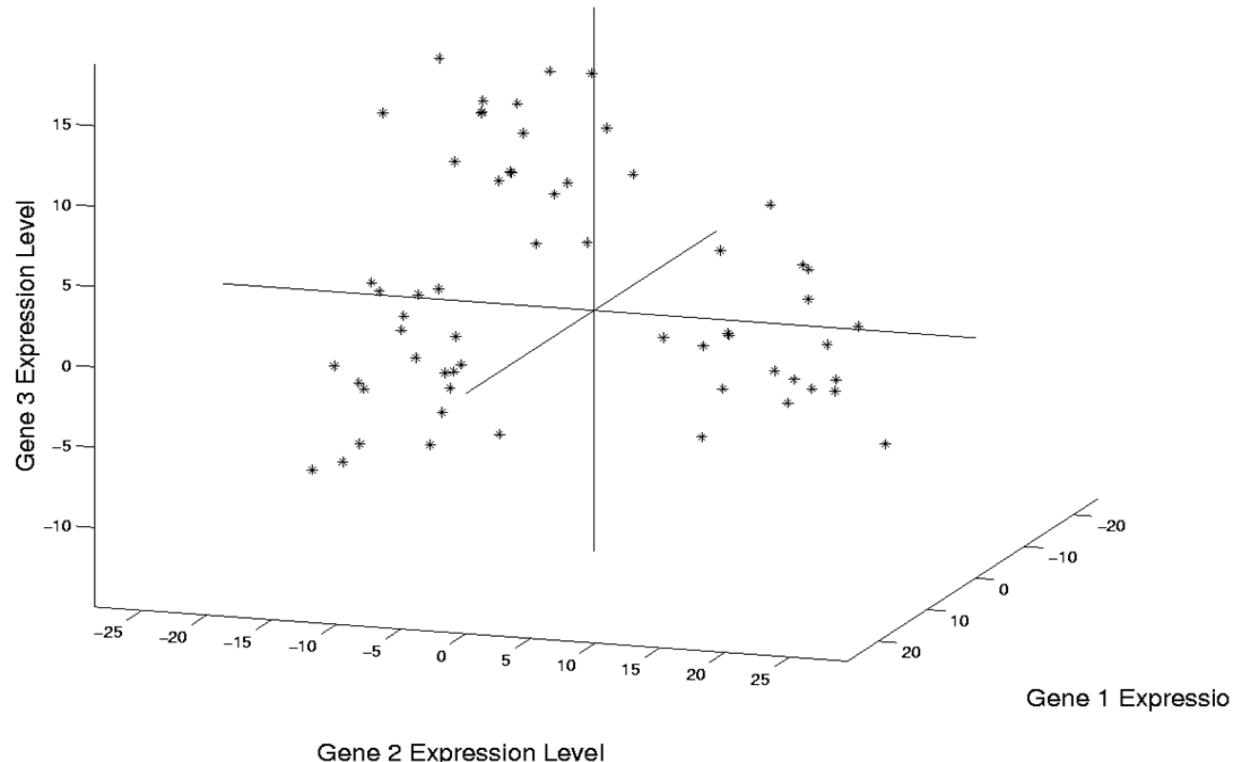
Example: Weather data



$n = 35$ and $p = 365$

Example

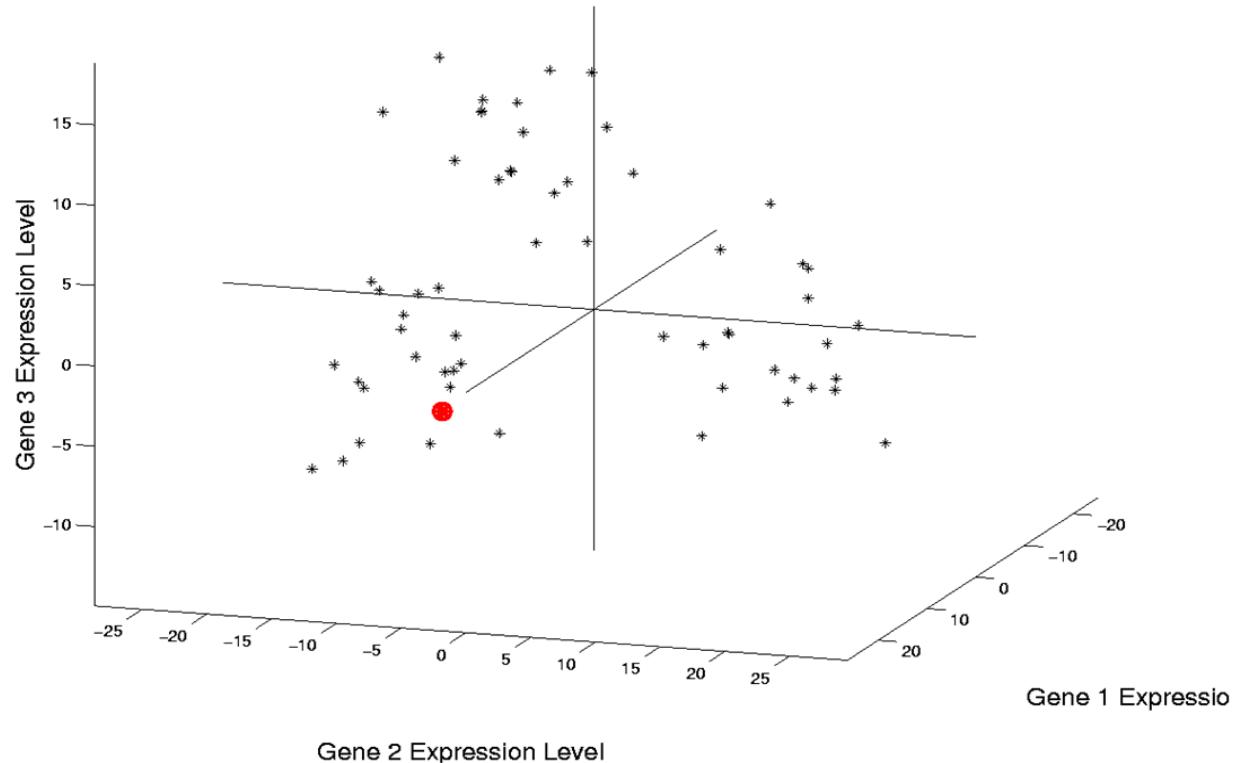
"Point Cloud View" of Gene Expression



$n = 60$ and $p = 3$

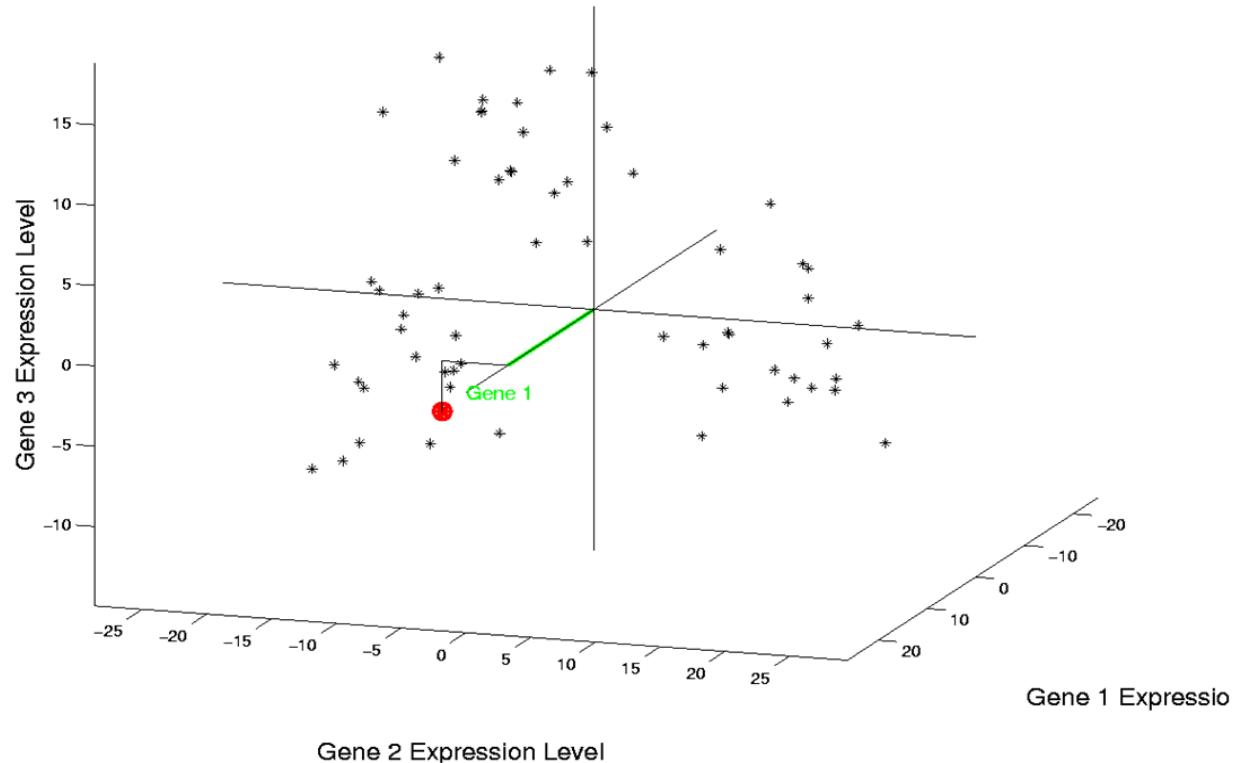
Example

"Point Cloud View" of Gene Expression



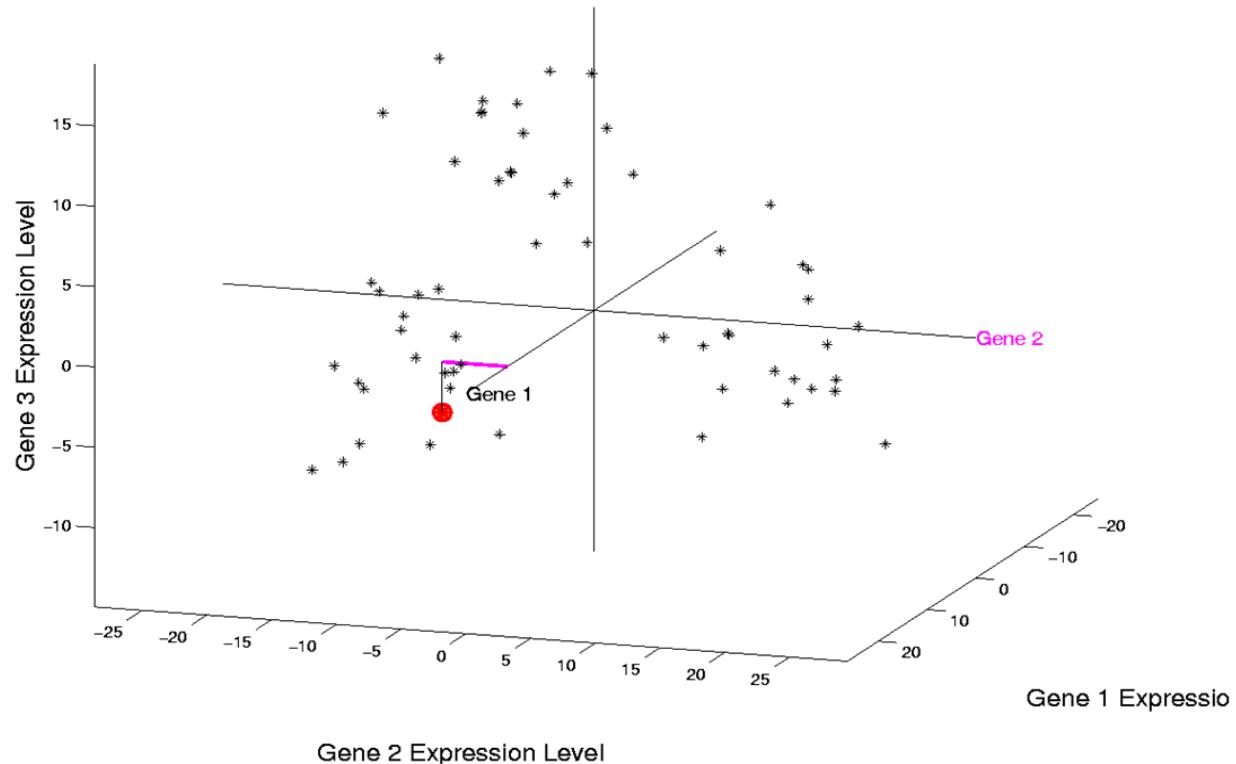
Example

"Point Cloud View" of Gene Expression



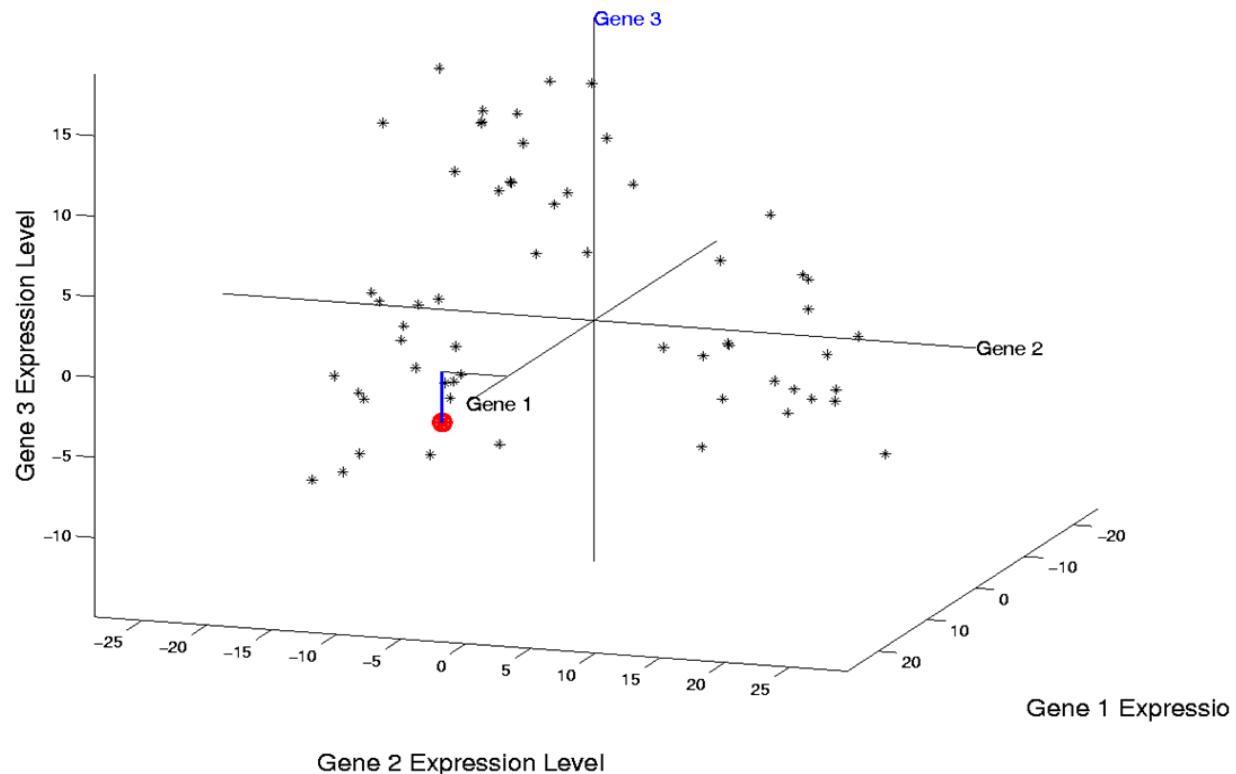
Example

"Point Cloud View" of Gene Expression



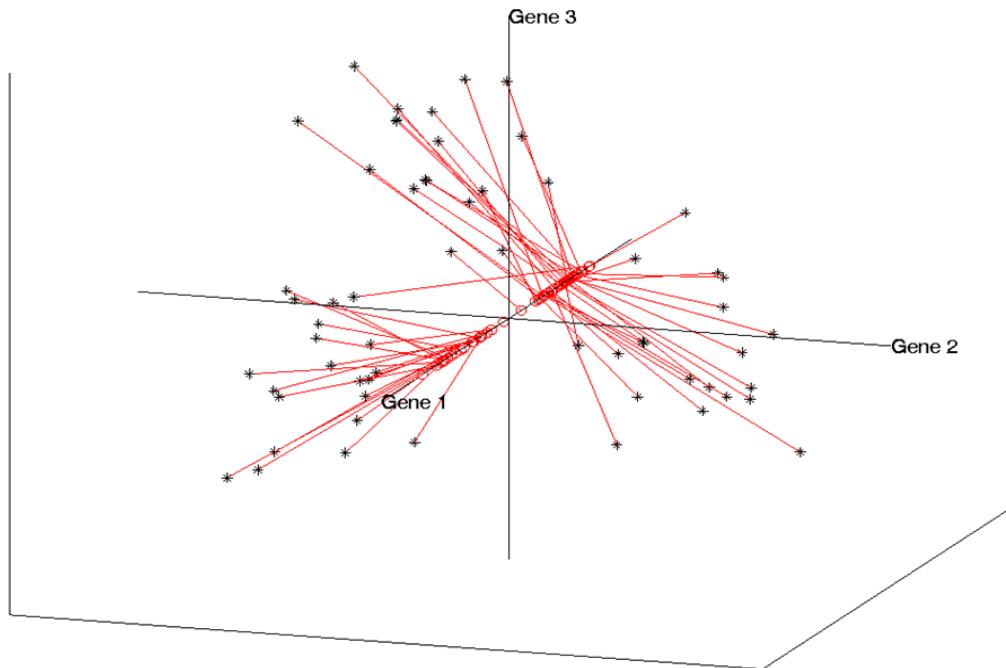
Example

"Point Cloud View" of Gene Expression

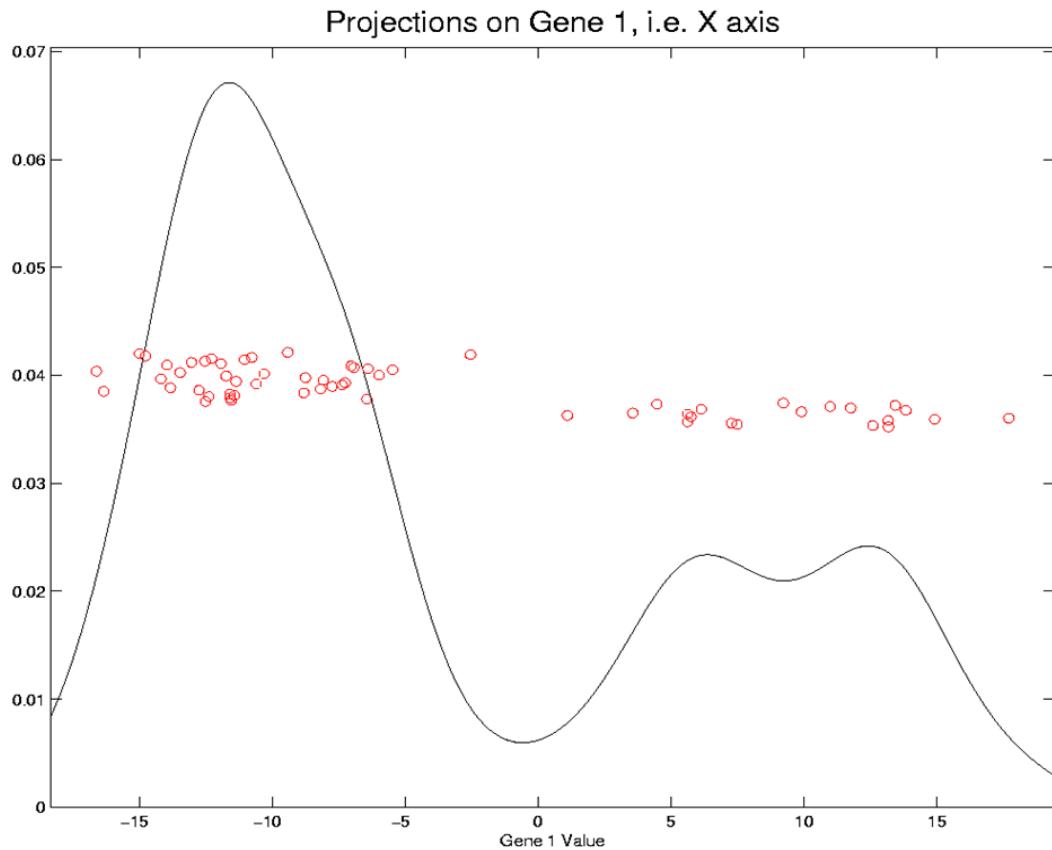


Example

Projections on Gene 1, i.e. X axis

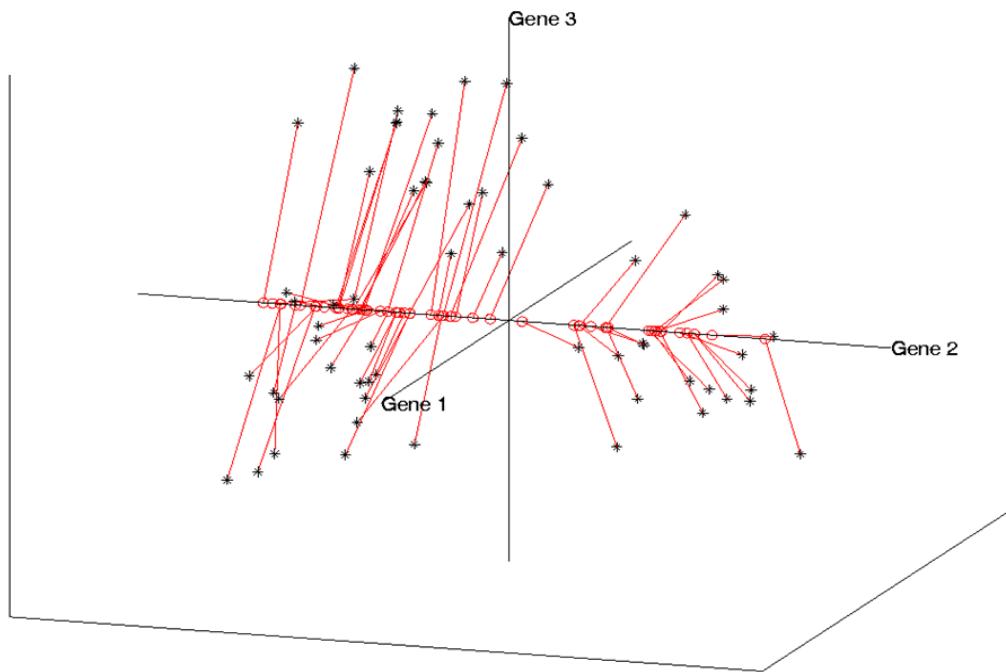


Example

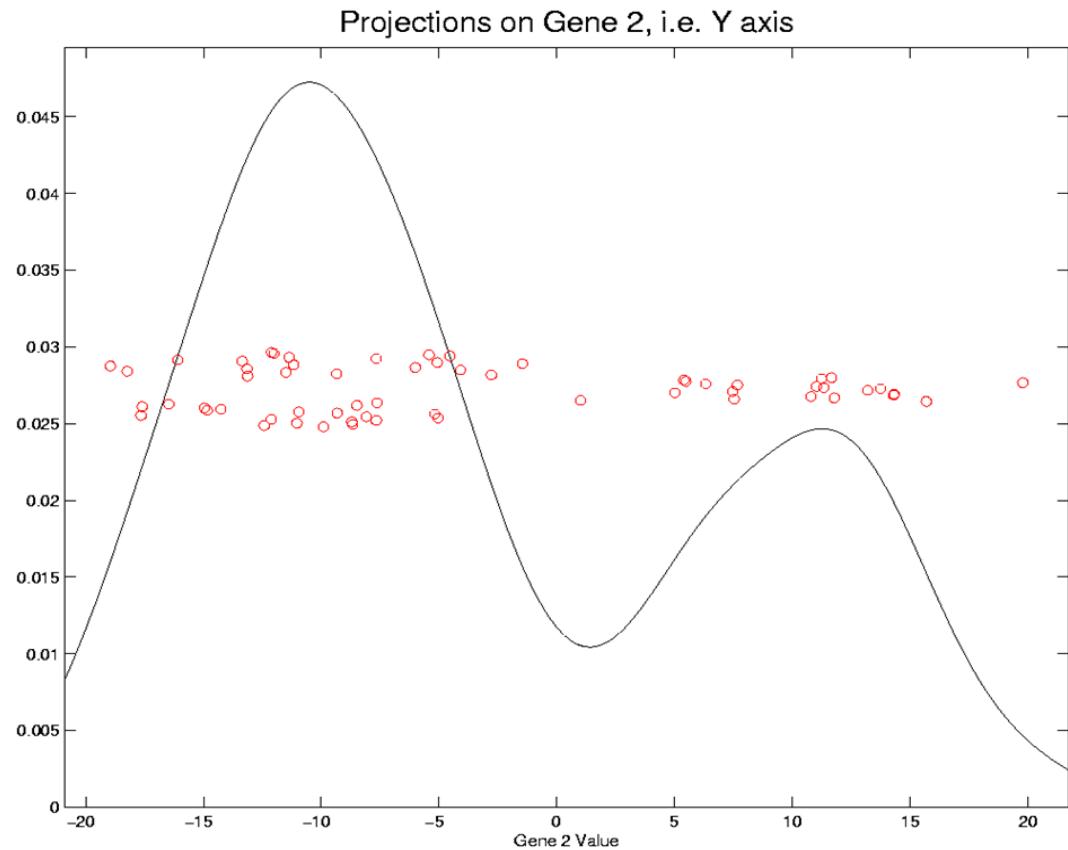


Example

Projections on Gene 2, i.e. Y axis

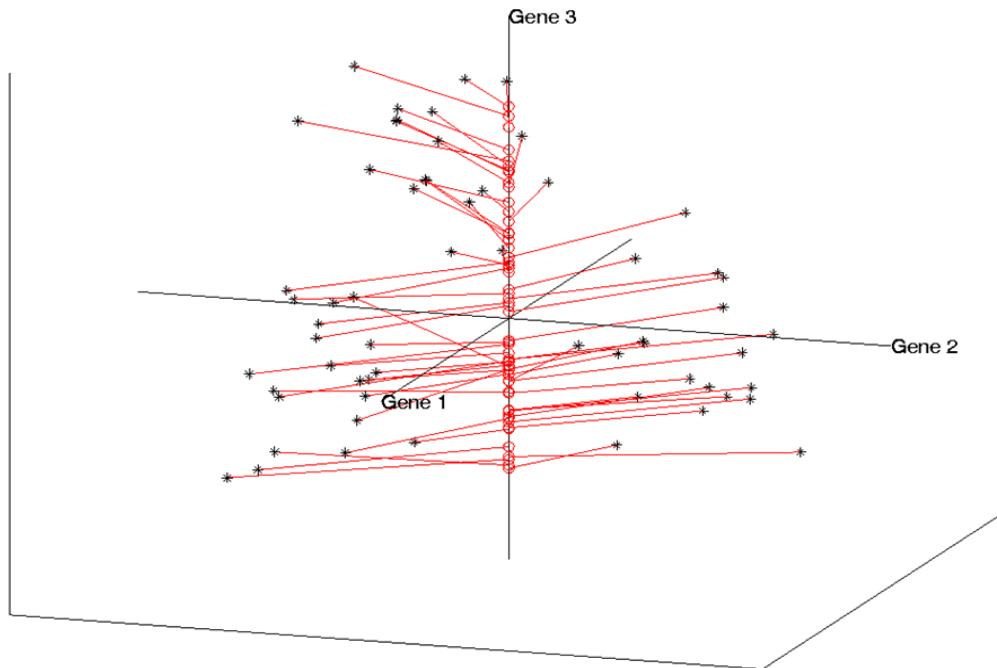


Example

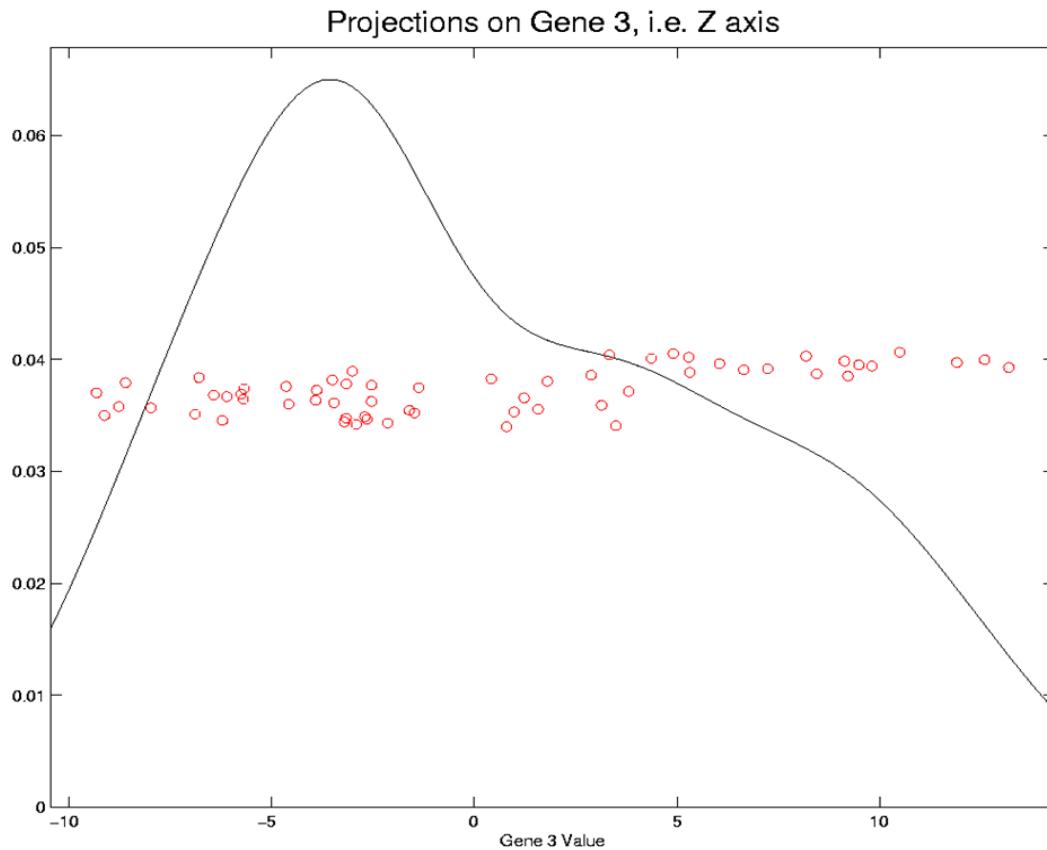


Example

Projections on Gene 3, i.e. Z axis

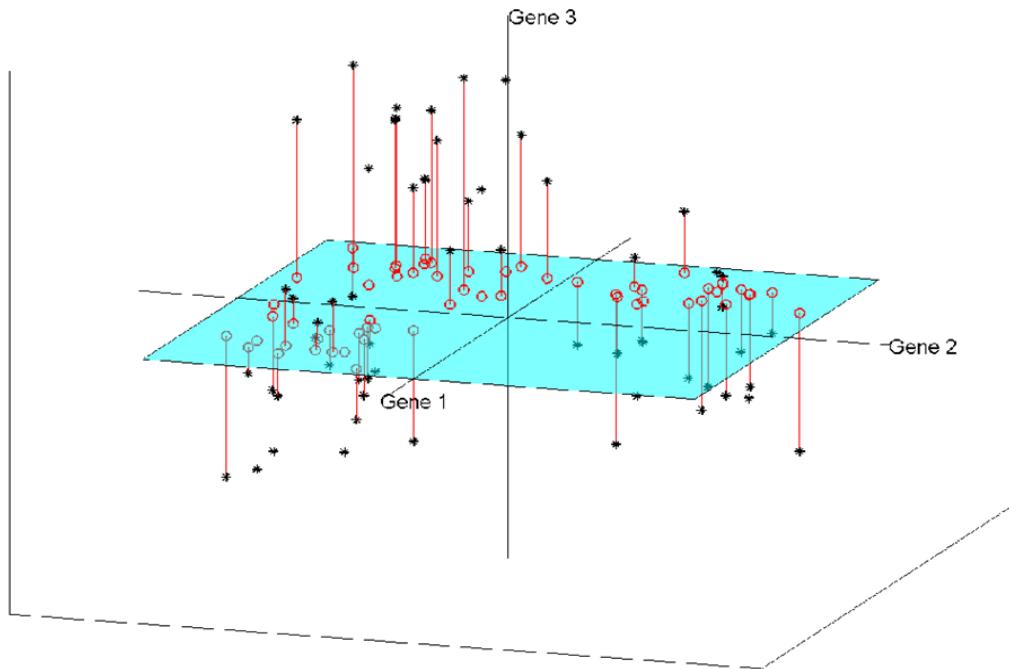


Example

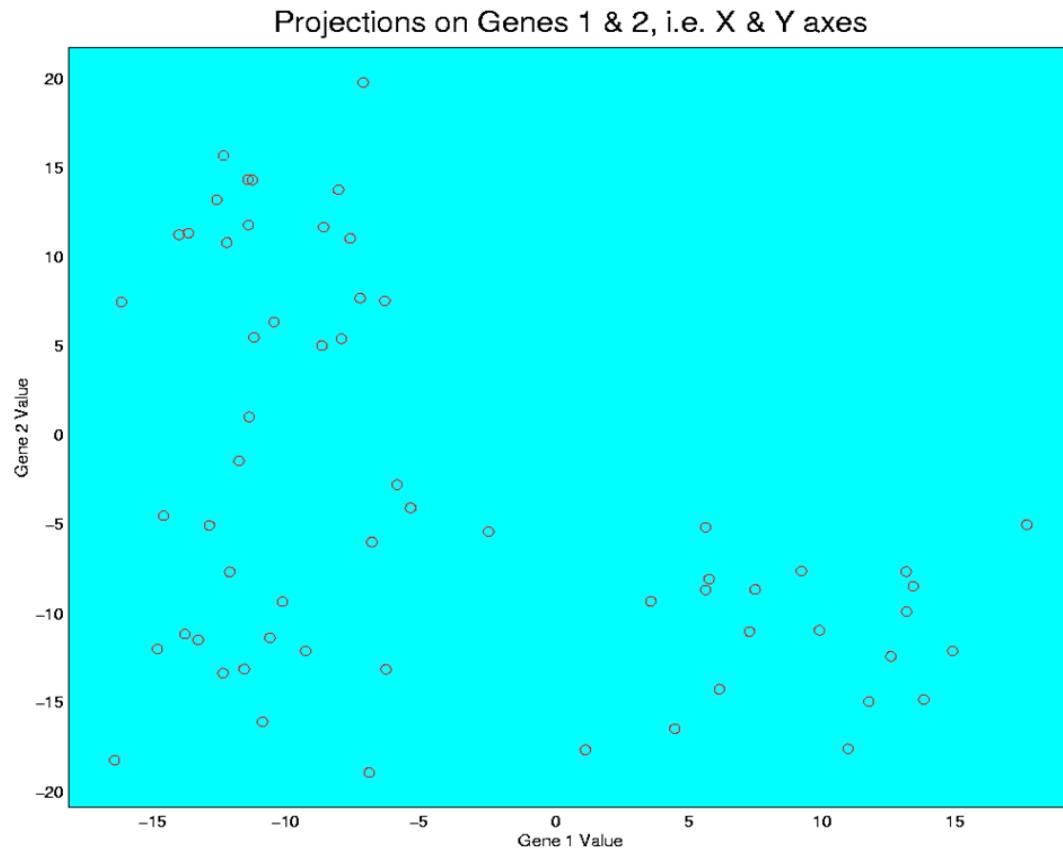


Example

Projections on Genes 1 & 2, i.e. X & Y axes

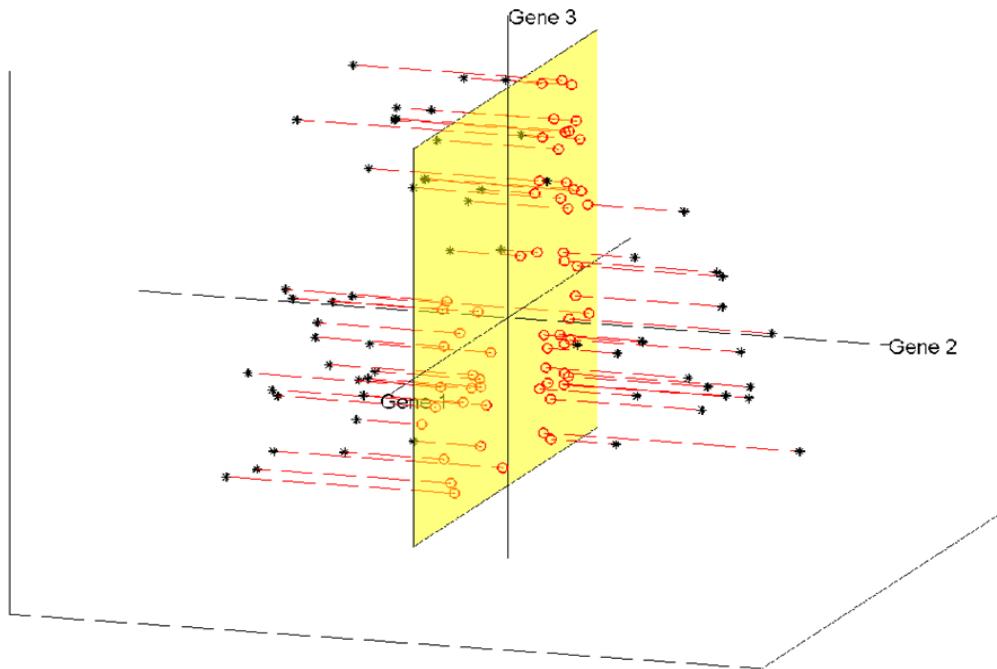


Example

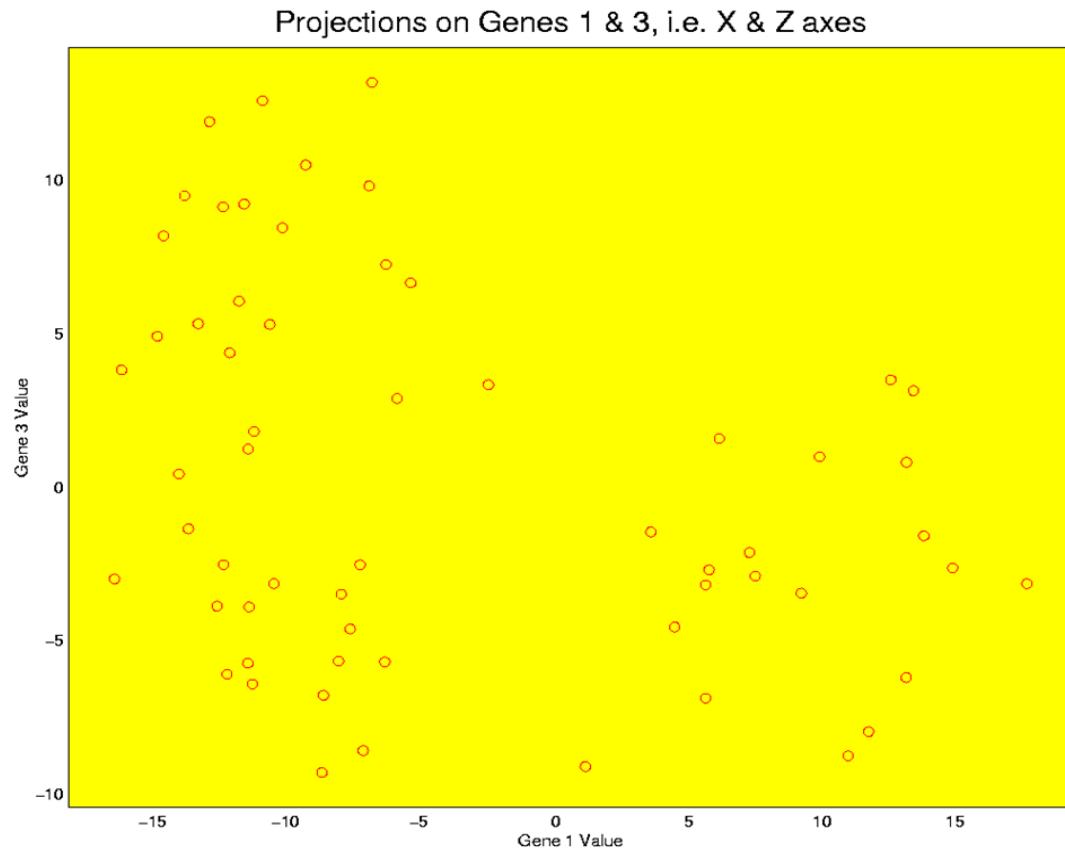


Example

Projections on Genes 1 & 3, i.e. X & Z axes

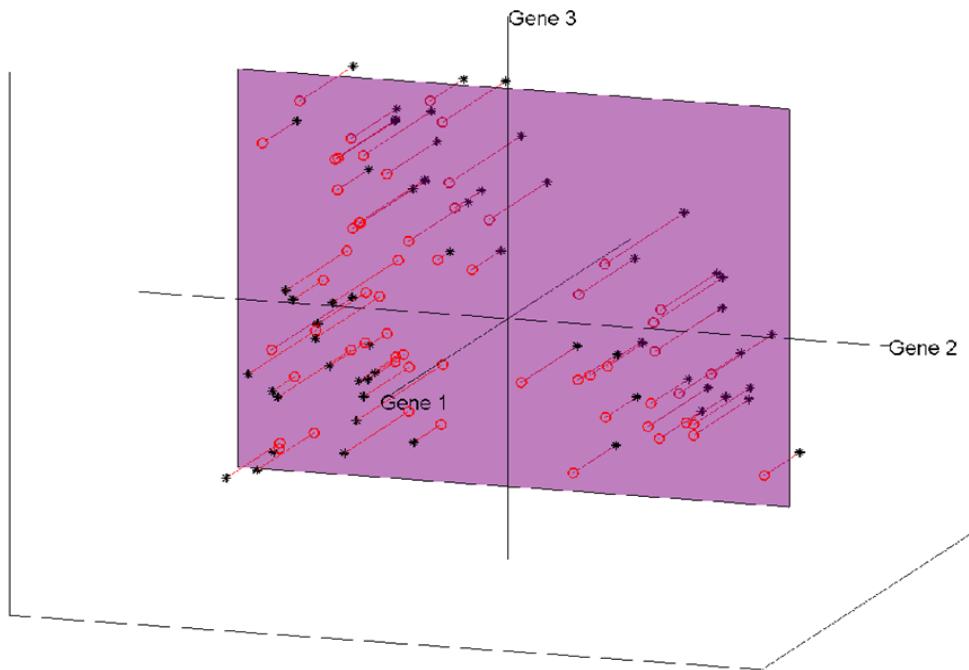


Example

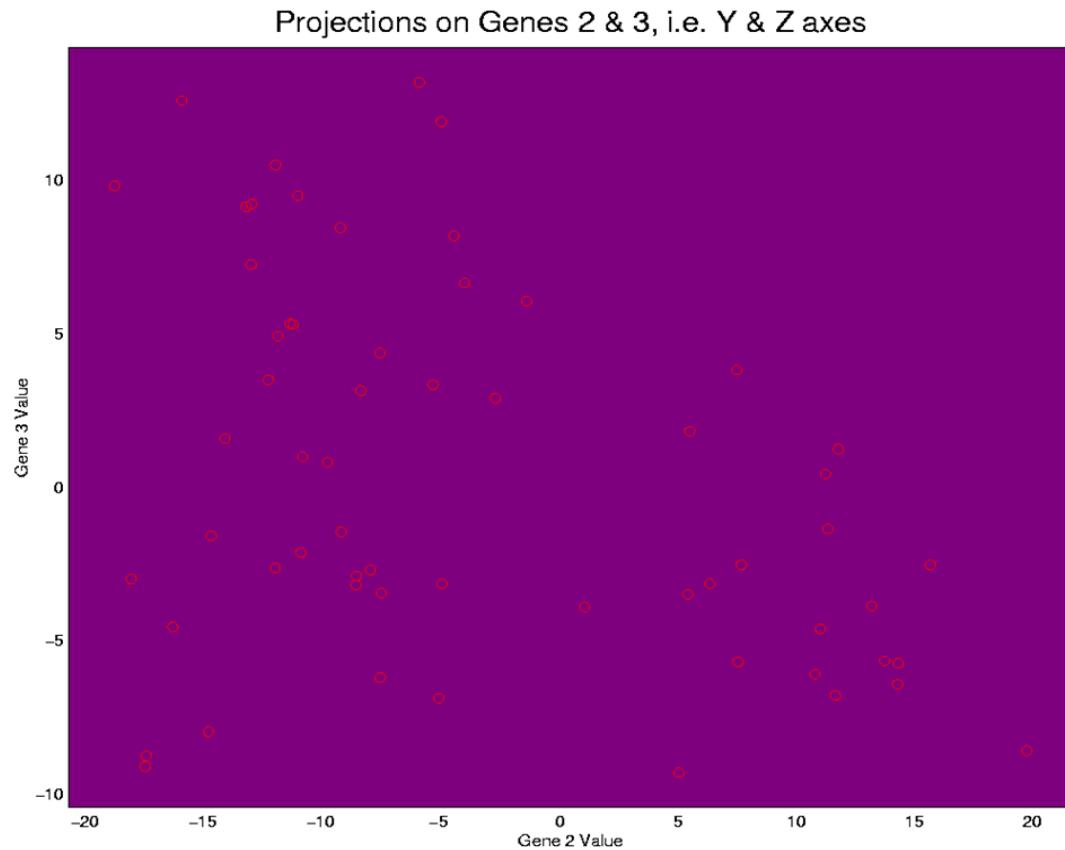


Example

Projections on Genes 2 & 3, i.e. Y & Z axes

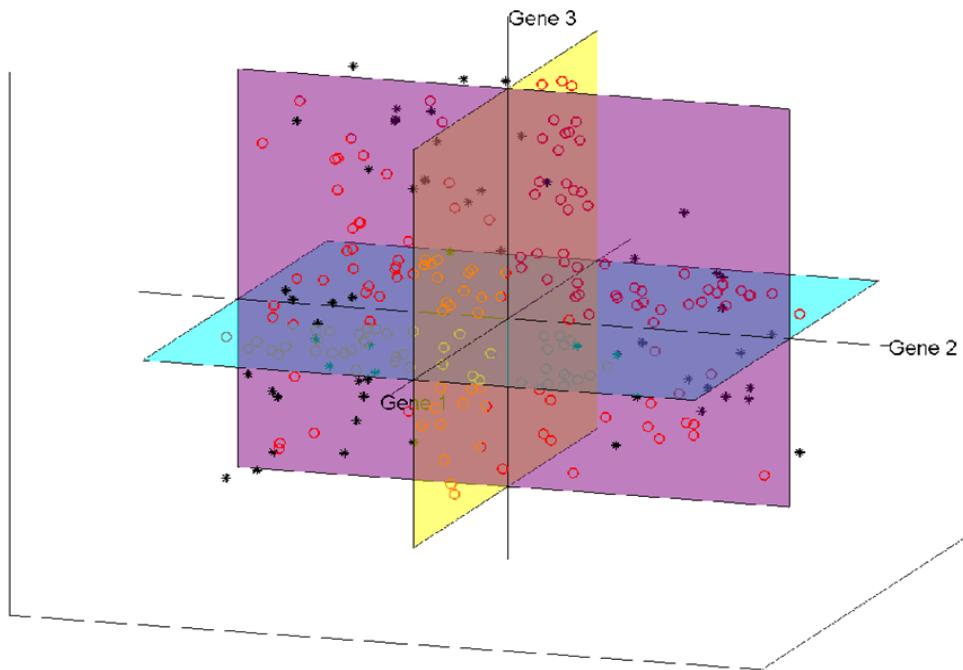


Example

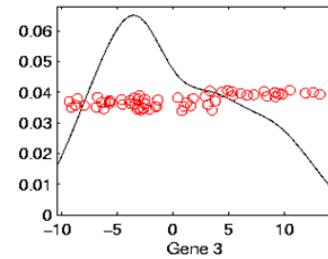
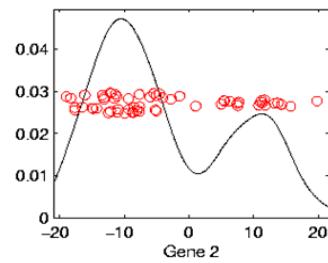
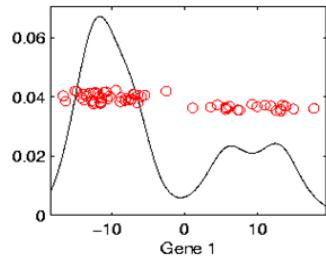


Example

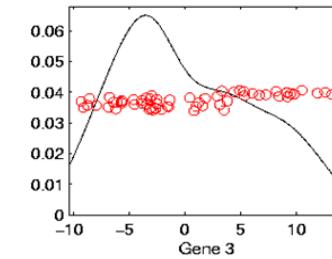
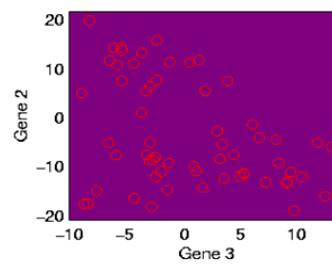
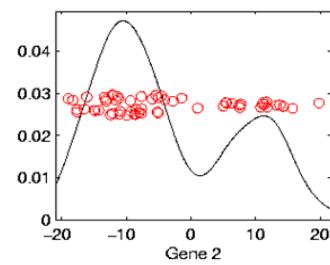
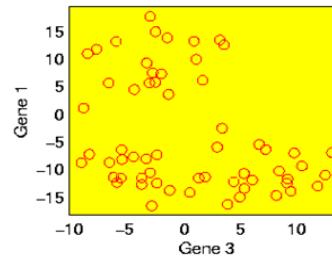
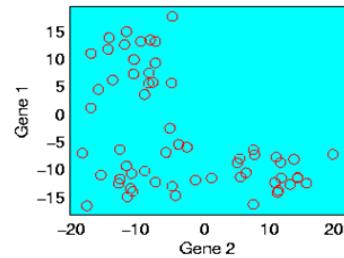
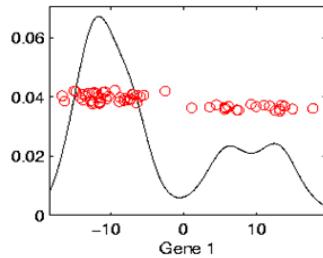
All Three 2d Projections



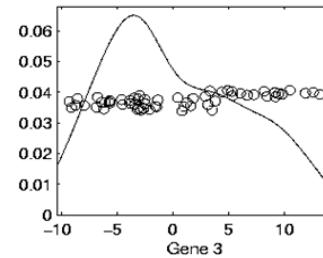
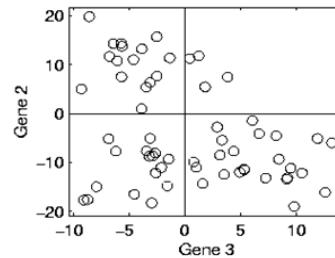
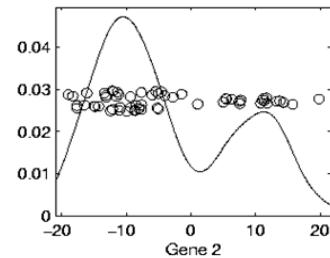
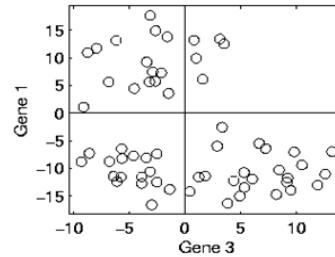
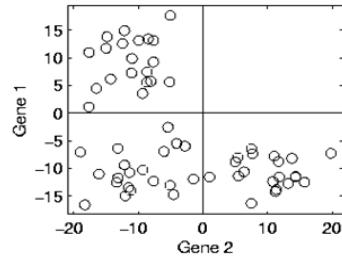
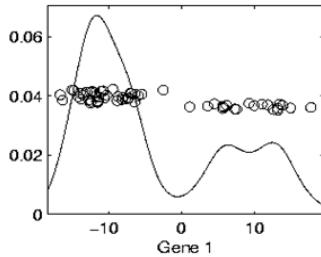
Example



Example



Example: Scatter plot matrix

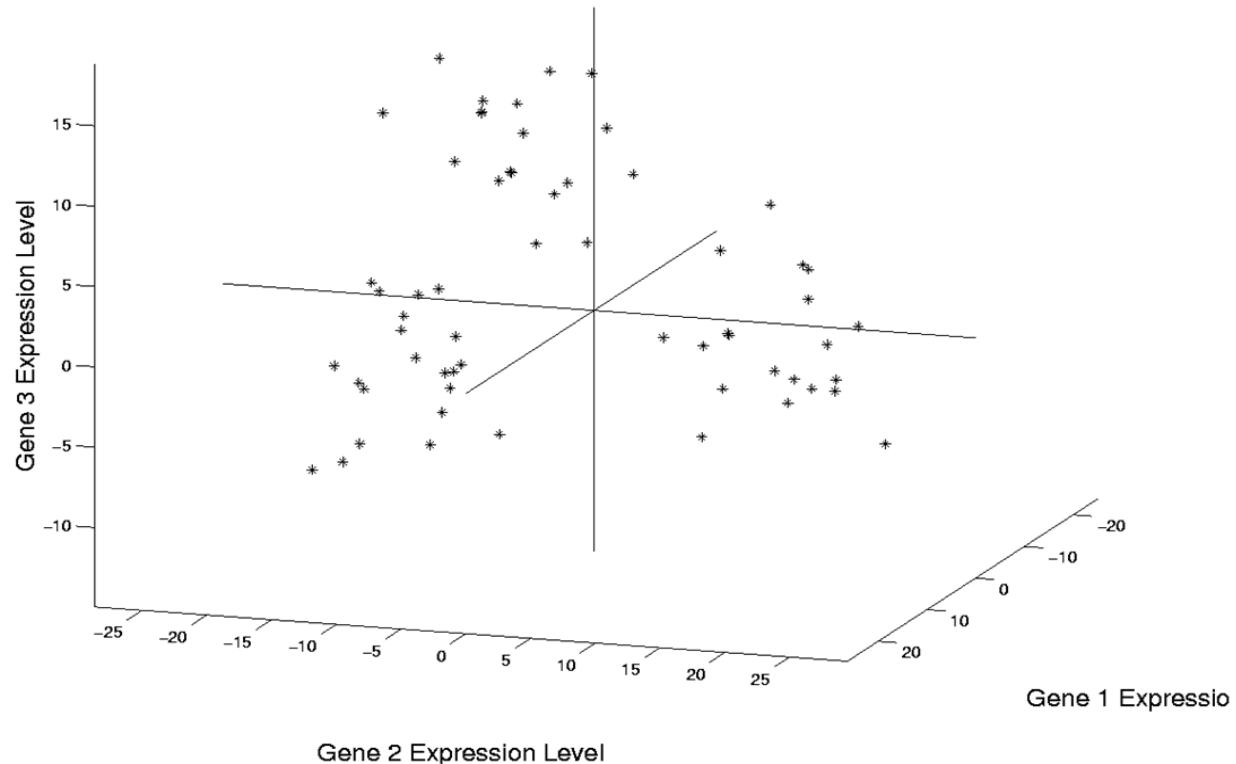


Principal Component Analysis

- Important set of directions: Principal Components
- Find directions of maximal (projected) variation
- Compute sequentially
- On orthogonal subspaces

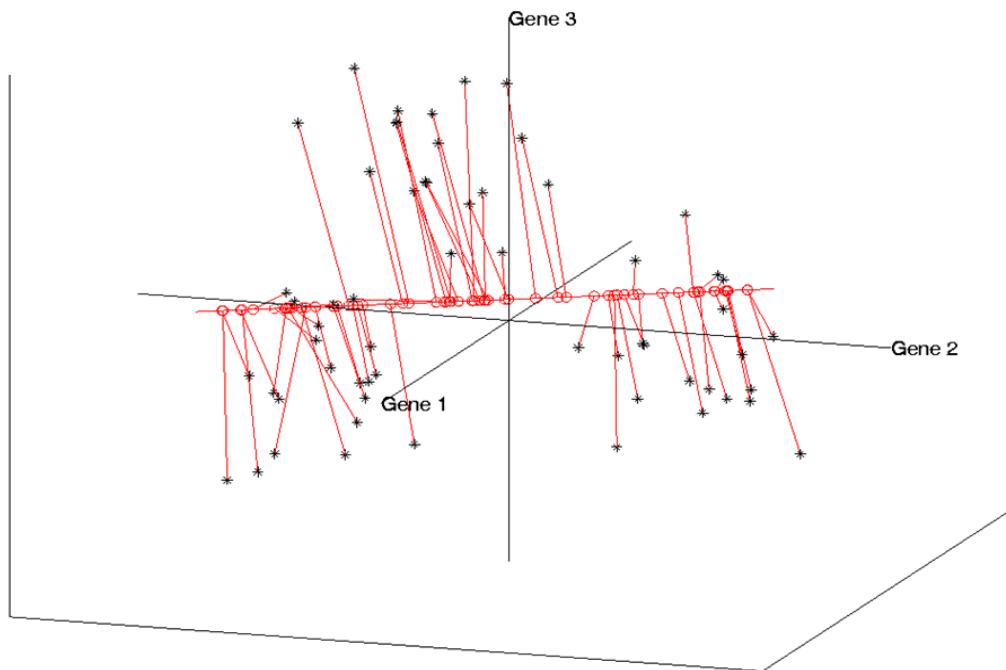
Example

"Point Cloud View" of Gene Expression

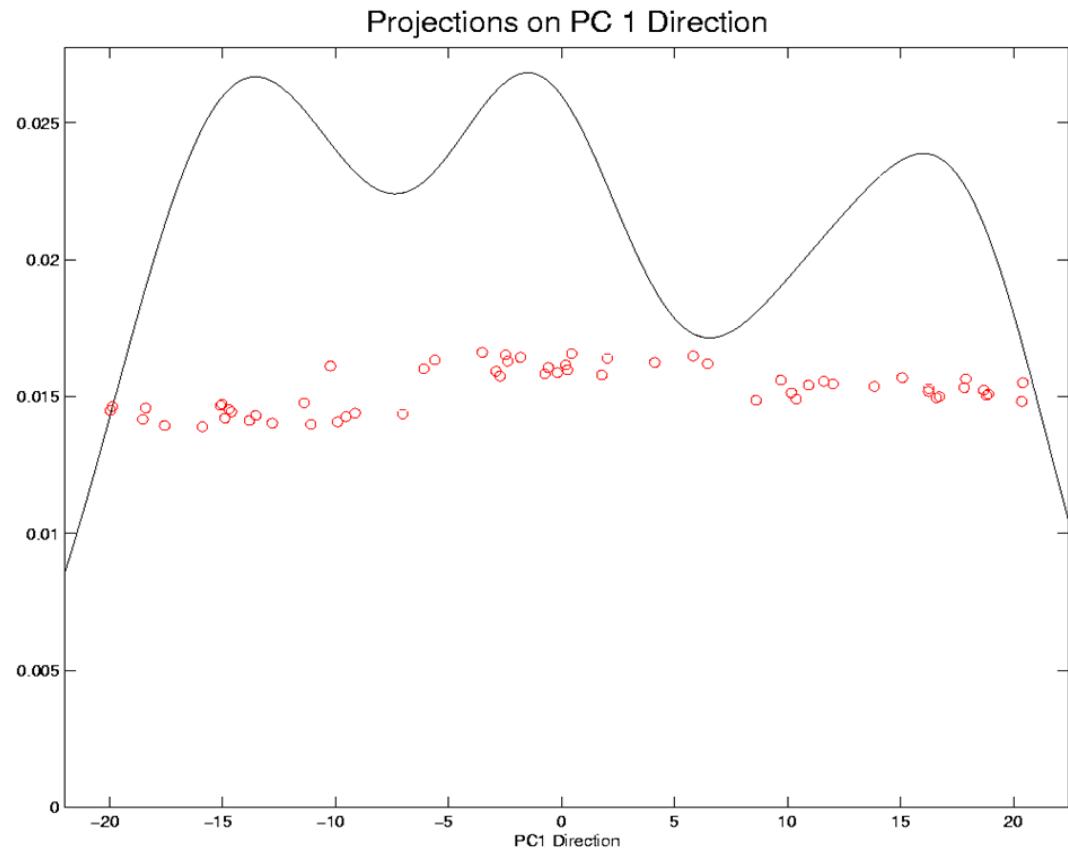


Example

Projections on PC 1 Direction

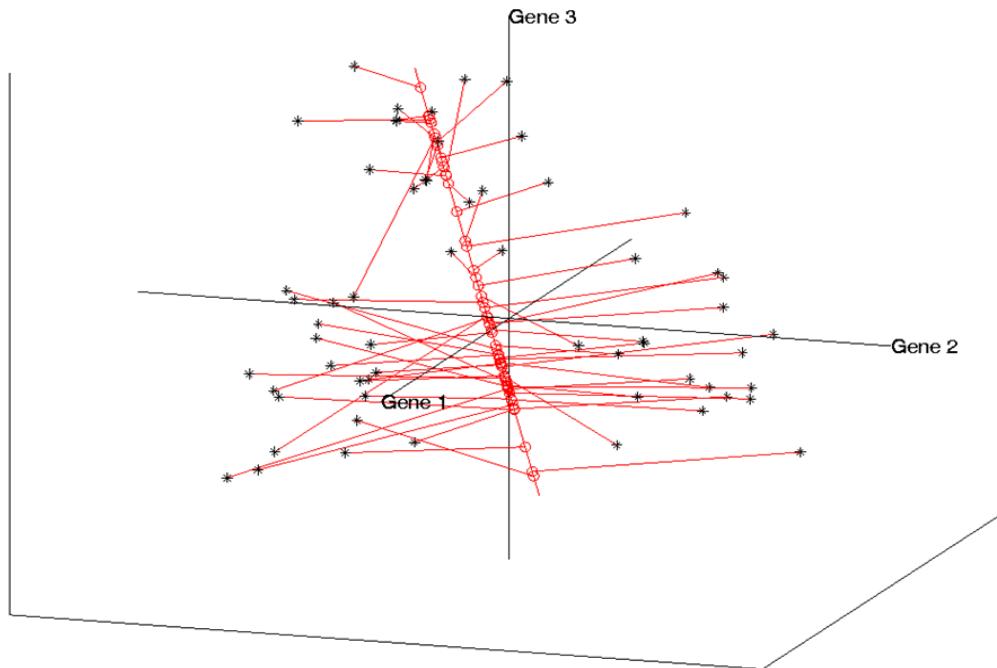


Example

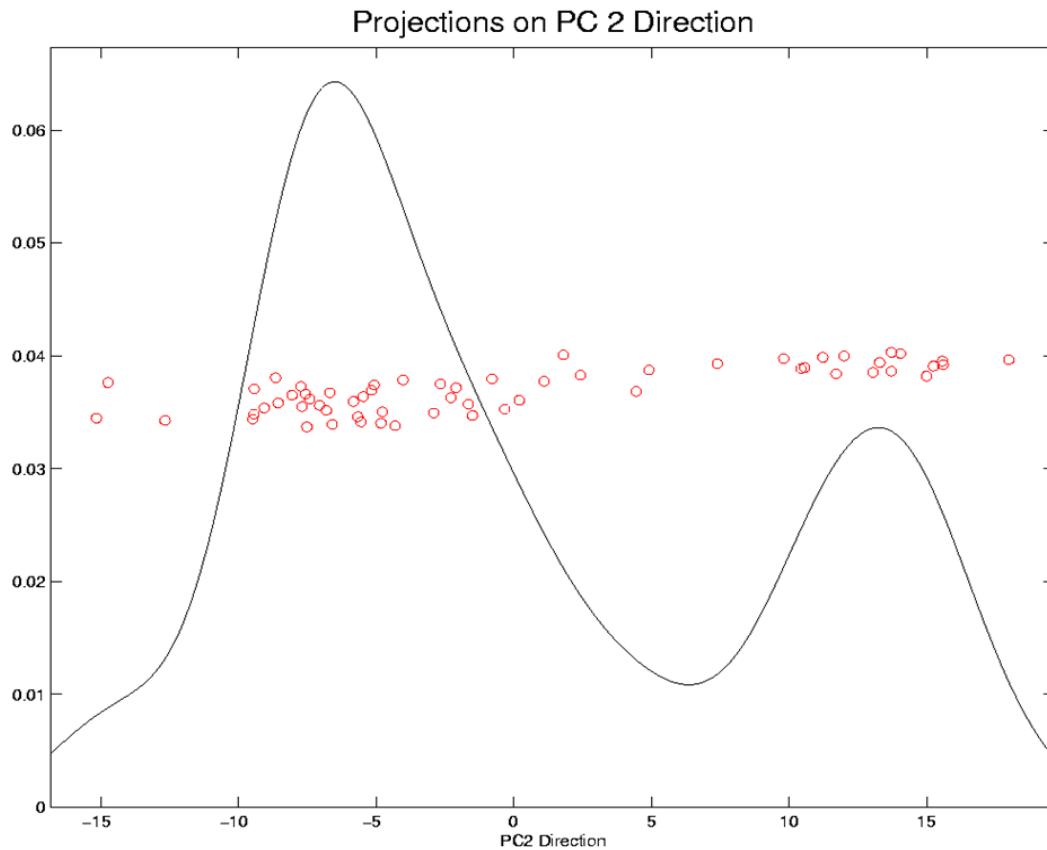


Example

Projections on PC 2 Direction

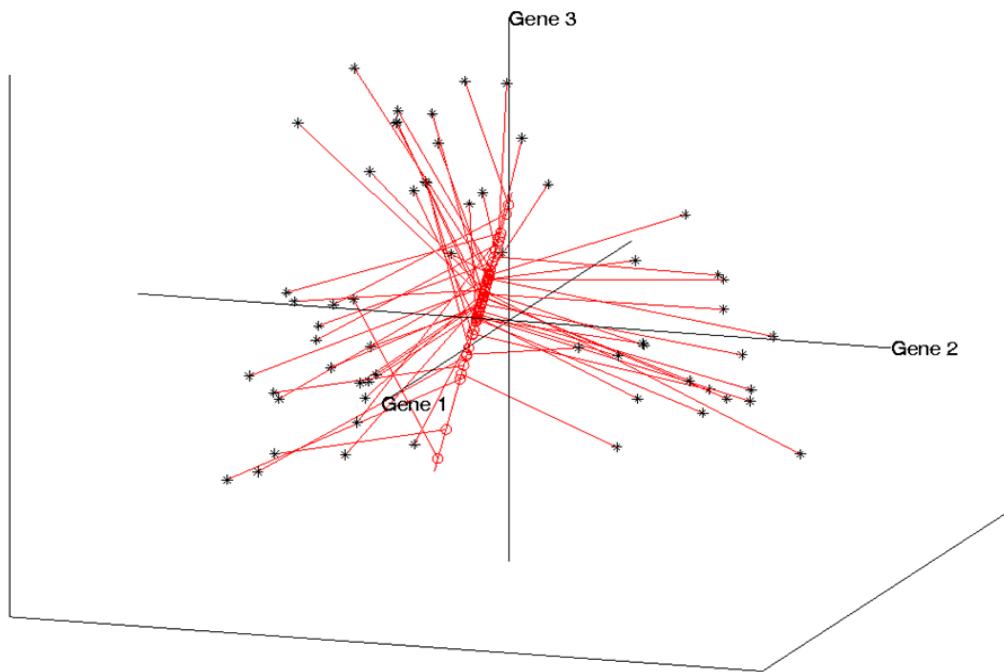


Example

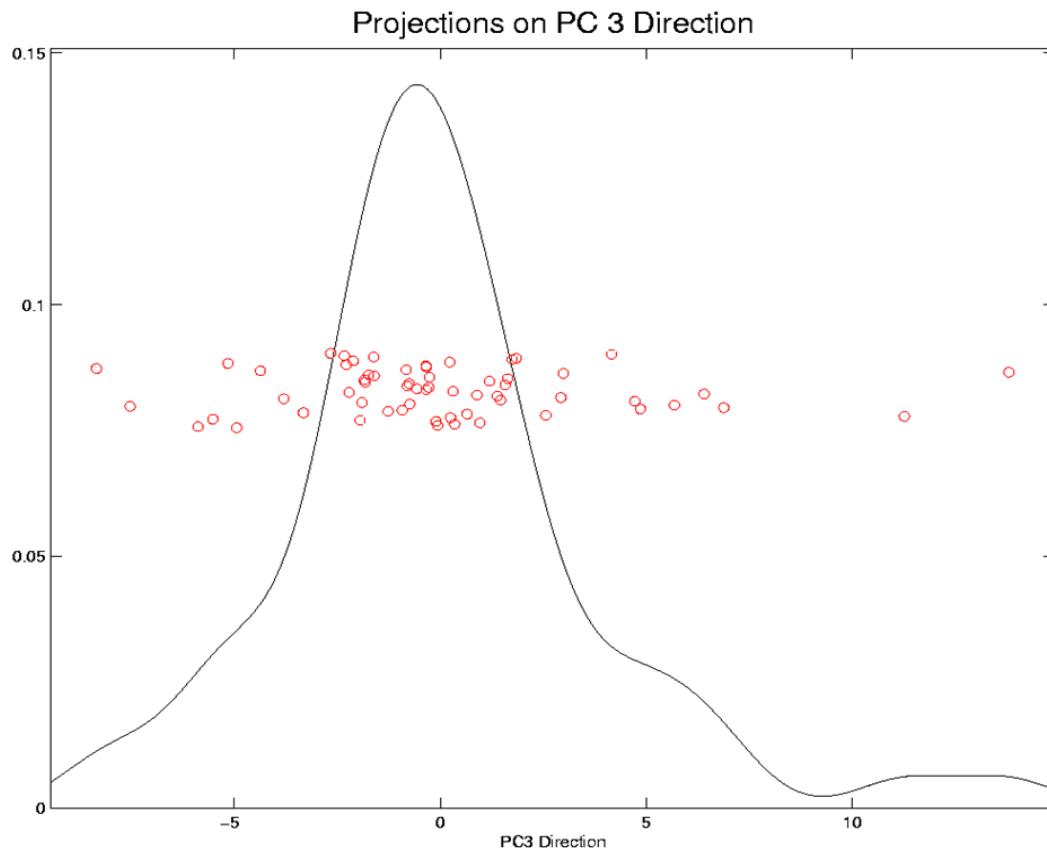


Example

Projections on PC 3 Direction

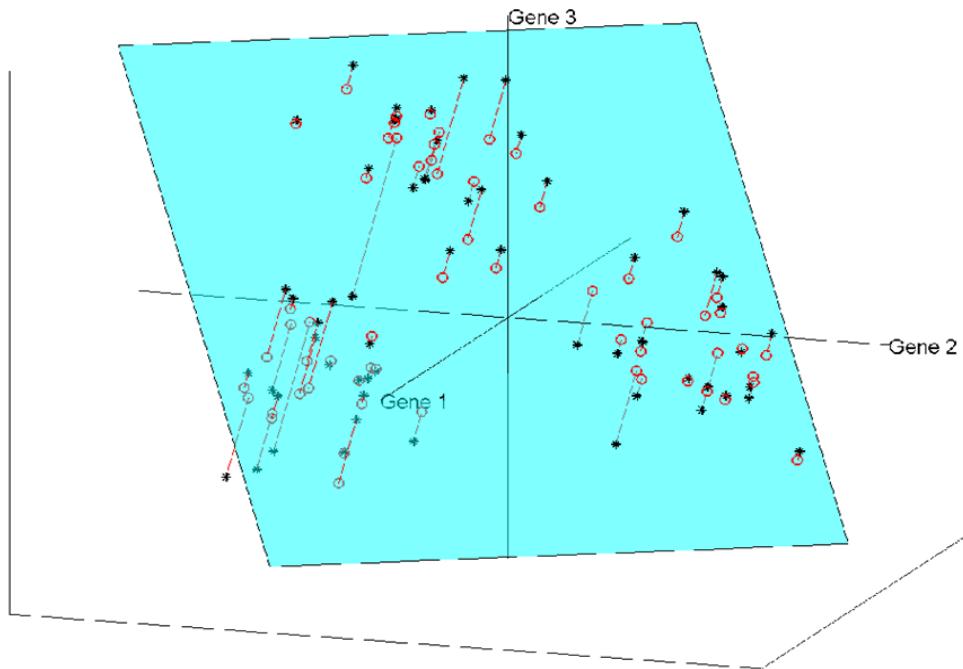


Example

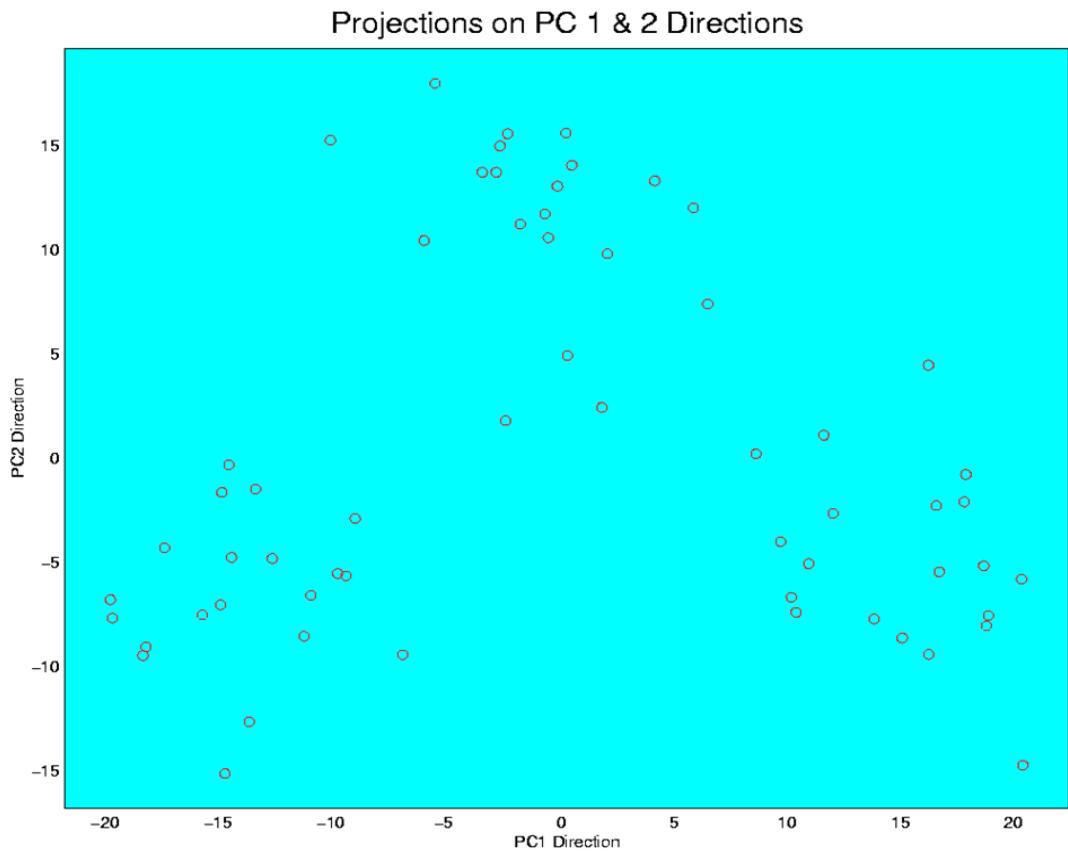


Example

Projections on PC 1 & PC 2 Directions

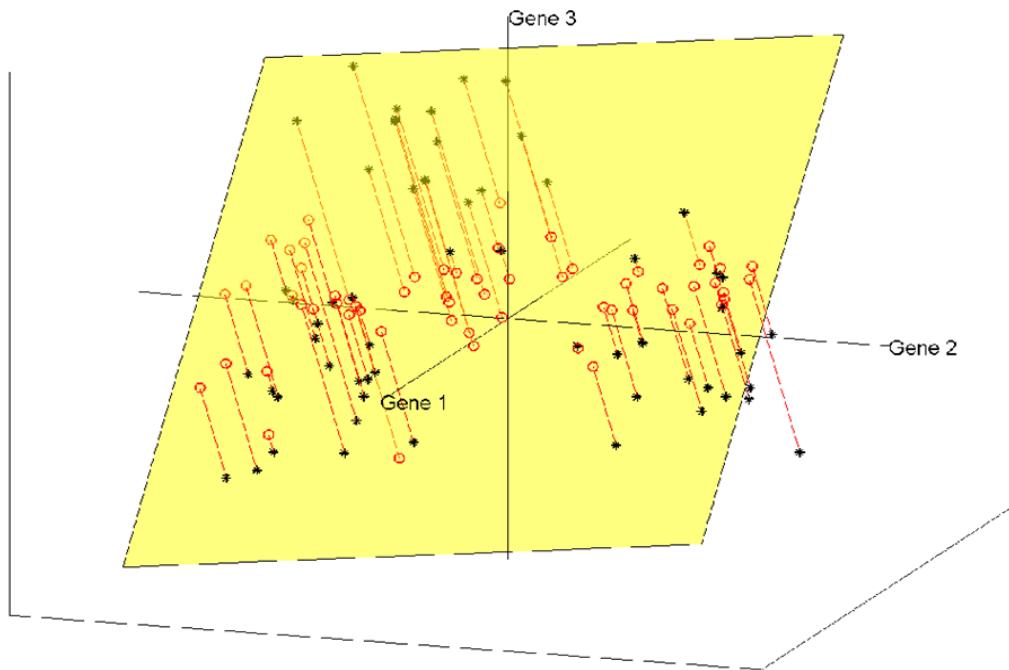


Example

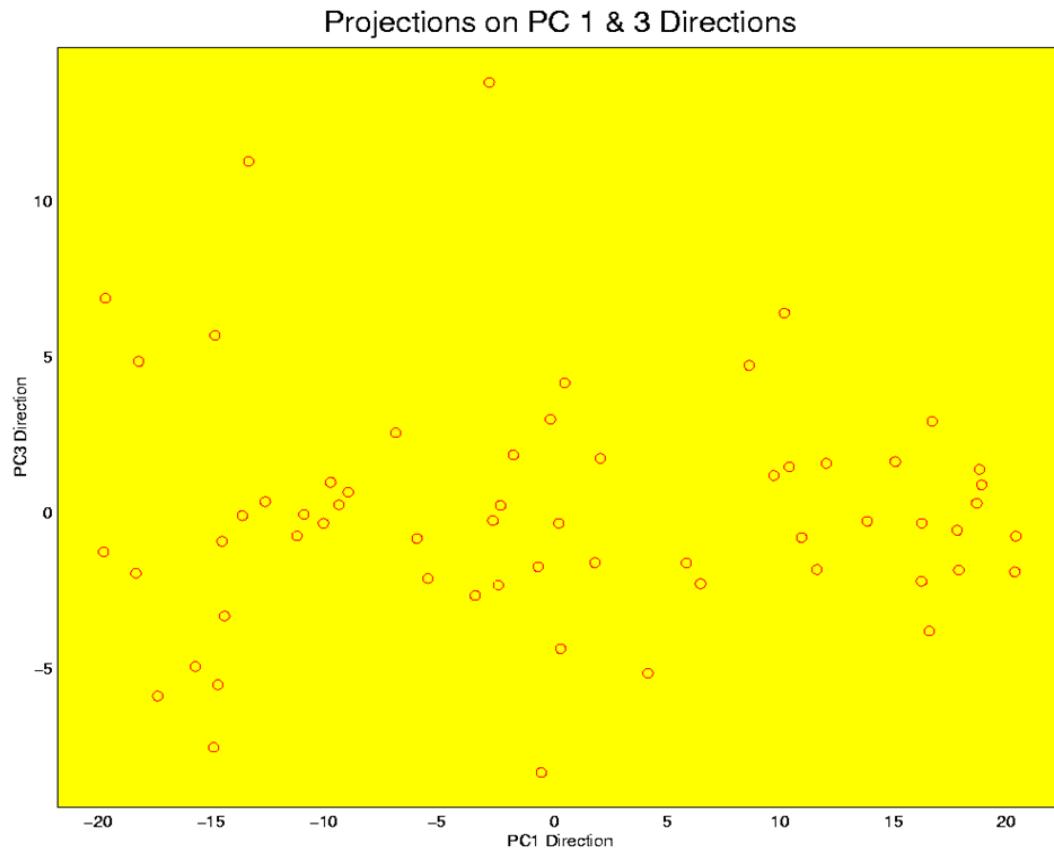


Example

Projections on PC 1 & PC 3 Directions

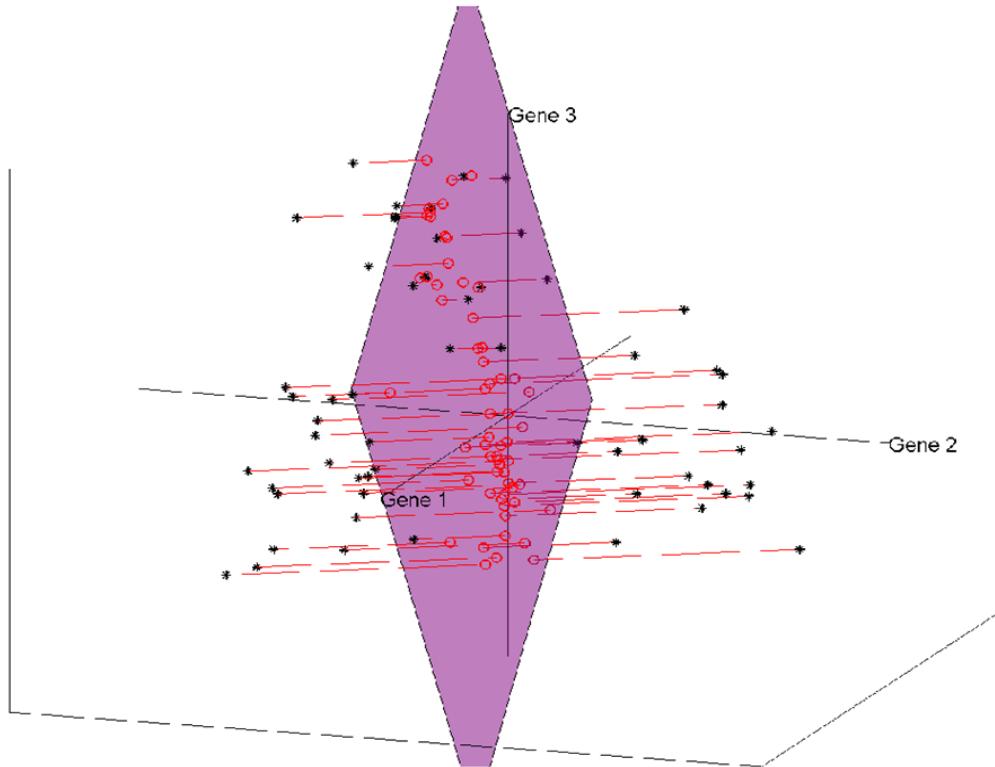


Example

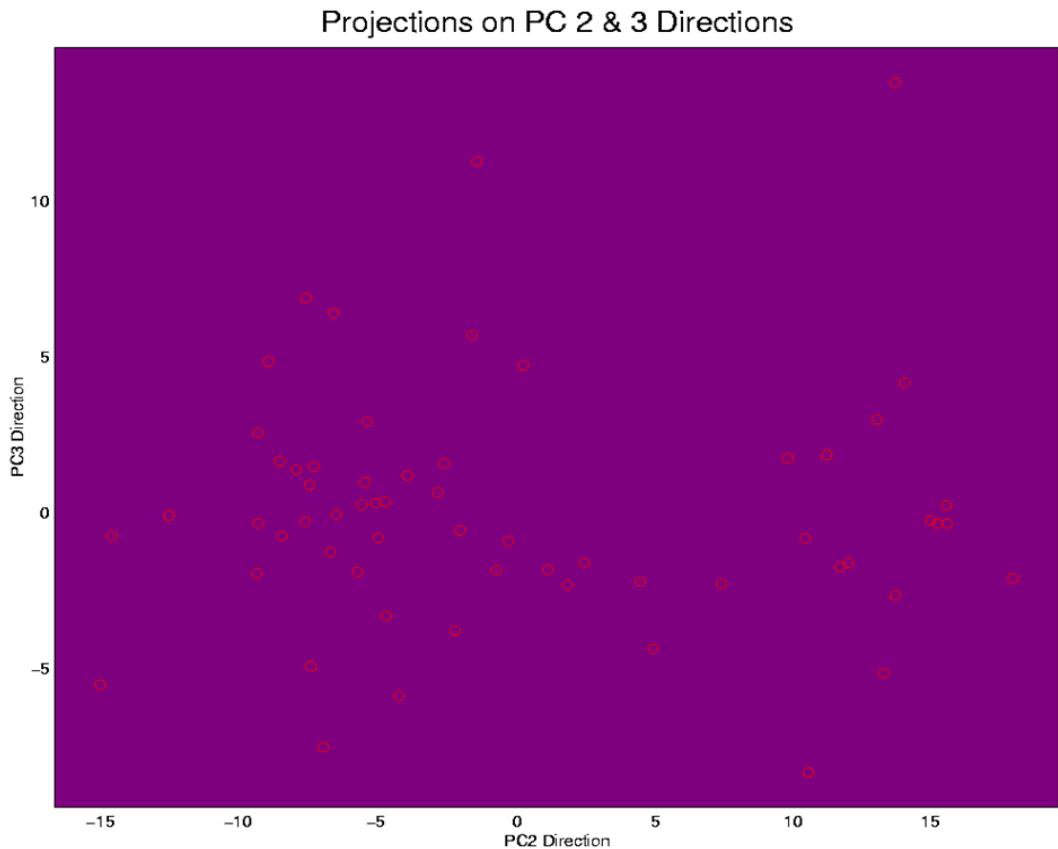


Example

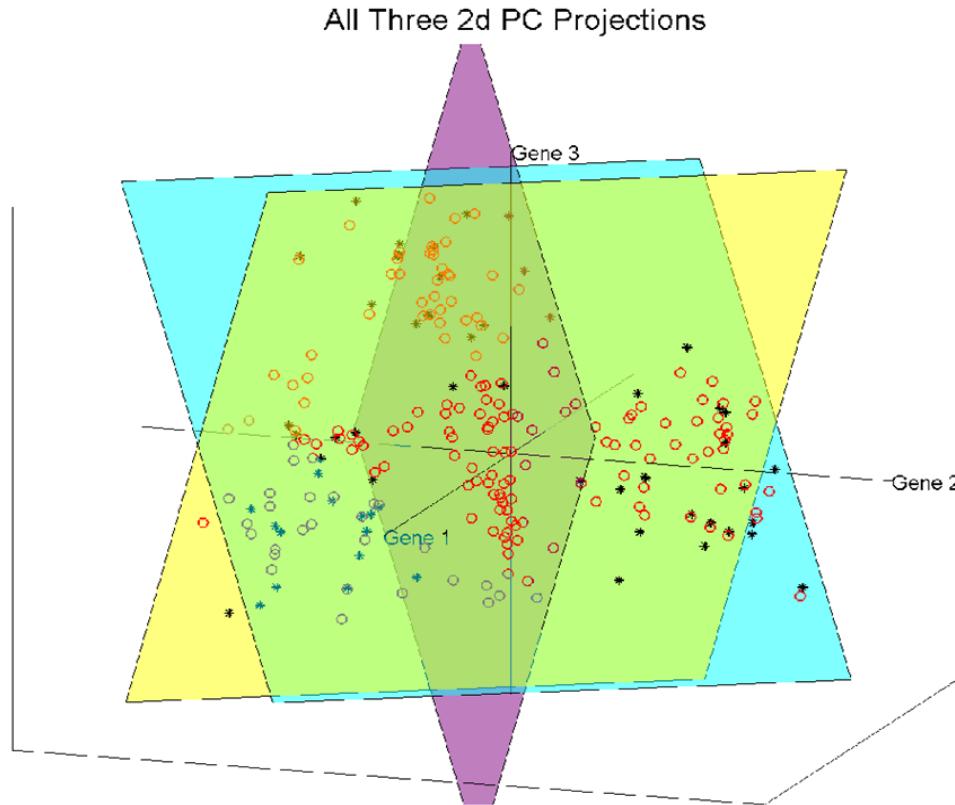
Projections on PC 2 & PC 3 Directions



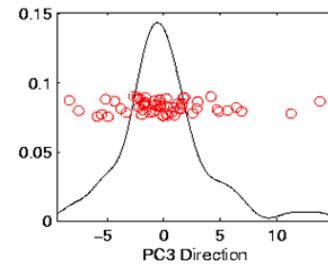
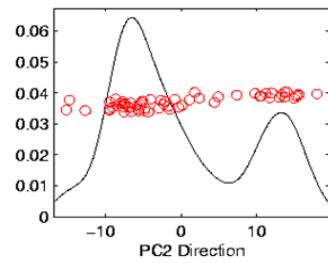
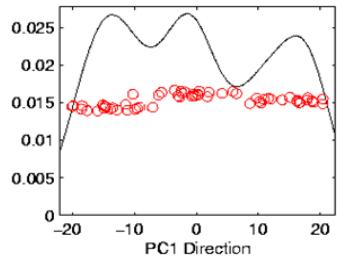
Example



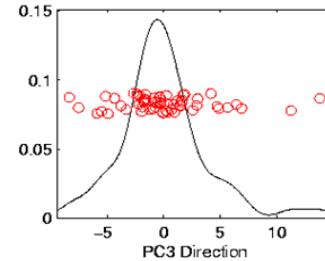
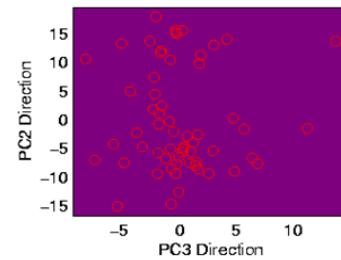
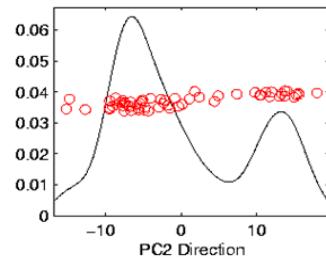
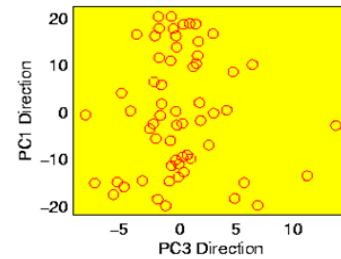
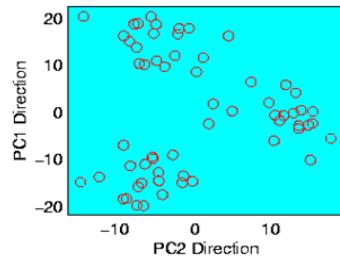
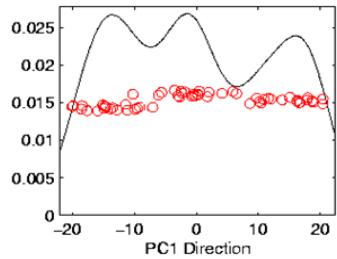
Example



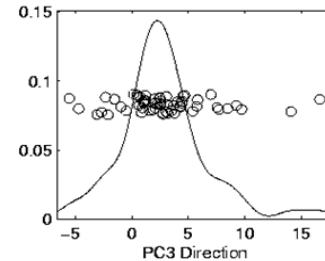
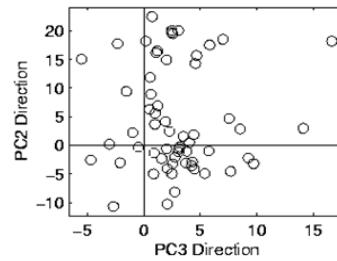
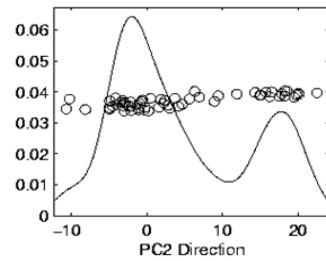
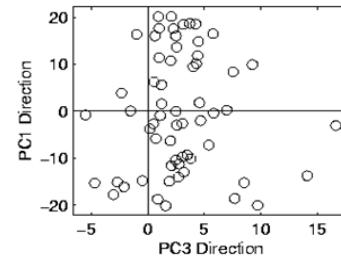
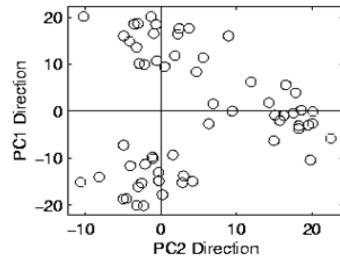
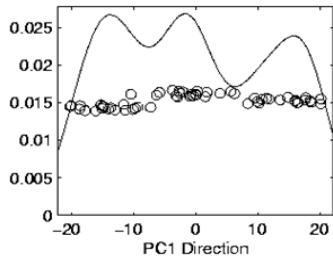
Example



Example

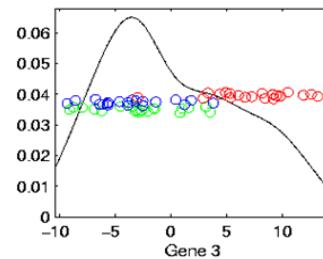
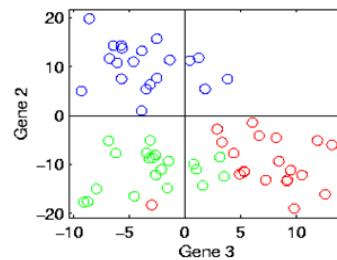
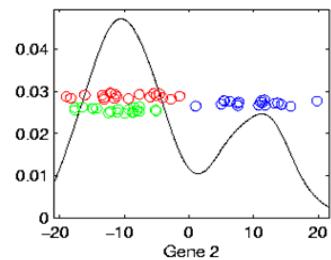
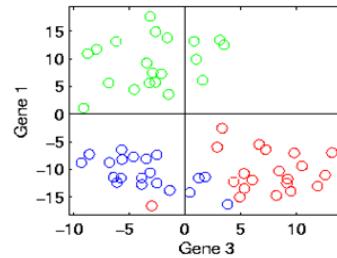
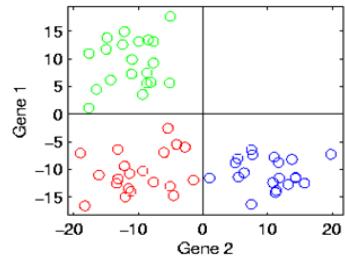
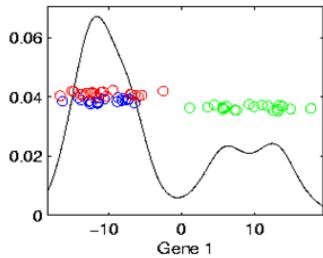


Example: PCA Scatter plot matrix

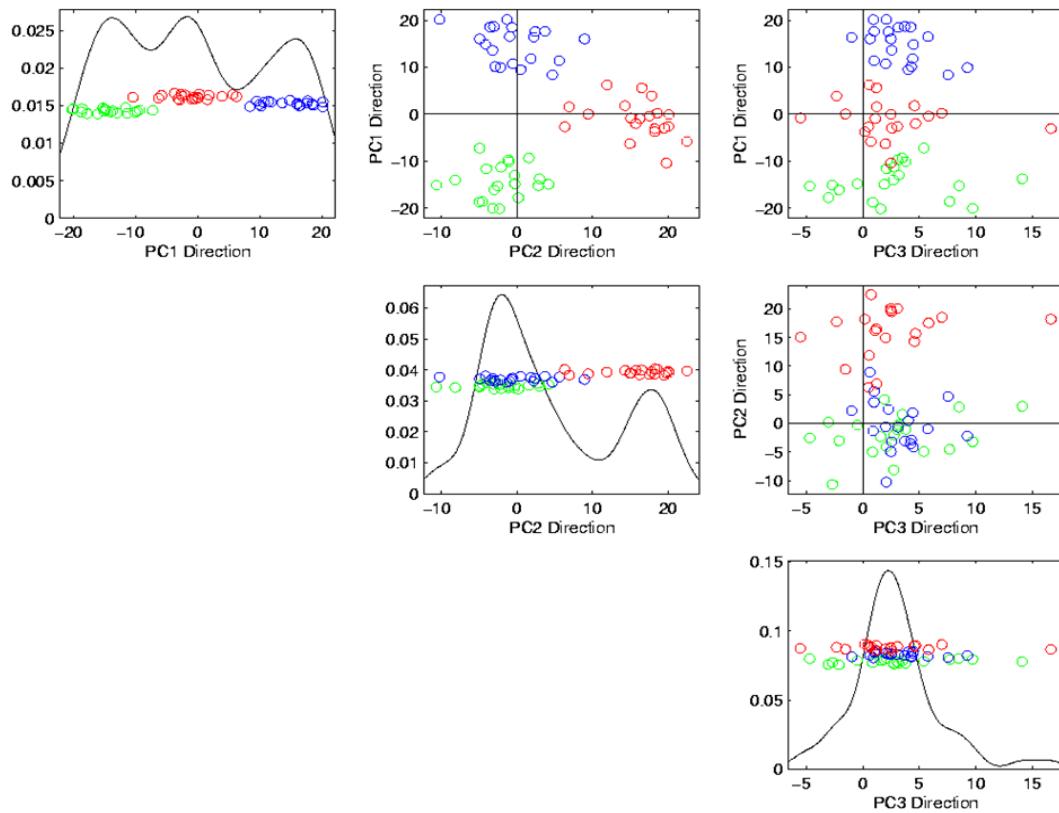


Now highlight 3 clusters

Example: Scatter plot matrix



Example: PCA Scatter plot matrix



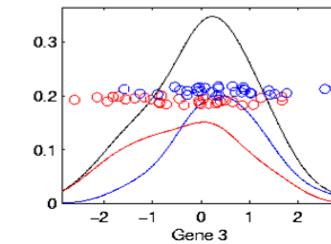
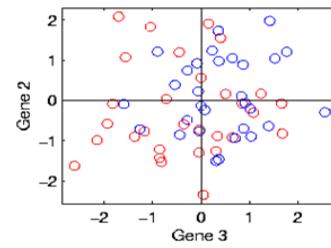
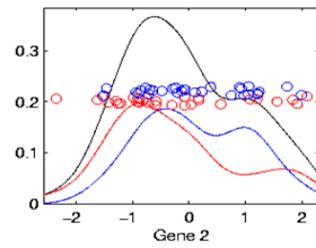
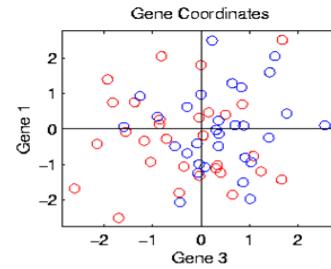
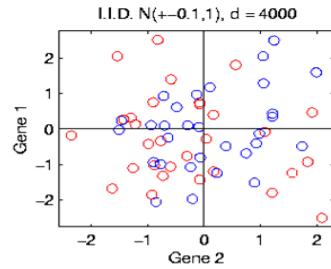
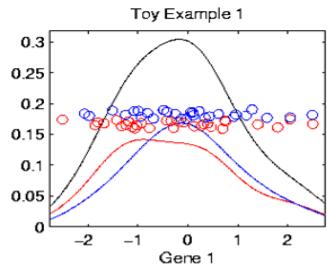
Comparison

- Gene by Gene view: Clusters appear in all 3 scatter plots, but never very separated
- PCA view: 1st PC shows three distinct clusters, better separated than in gene view

Another Example

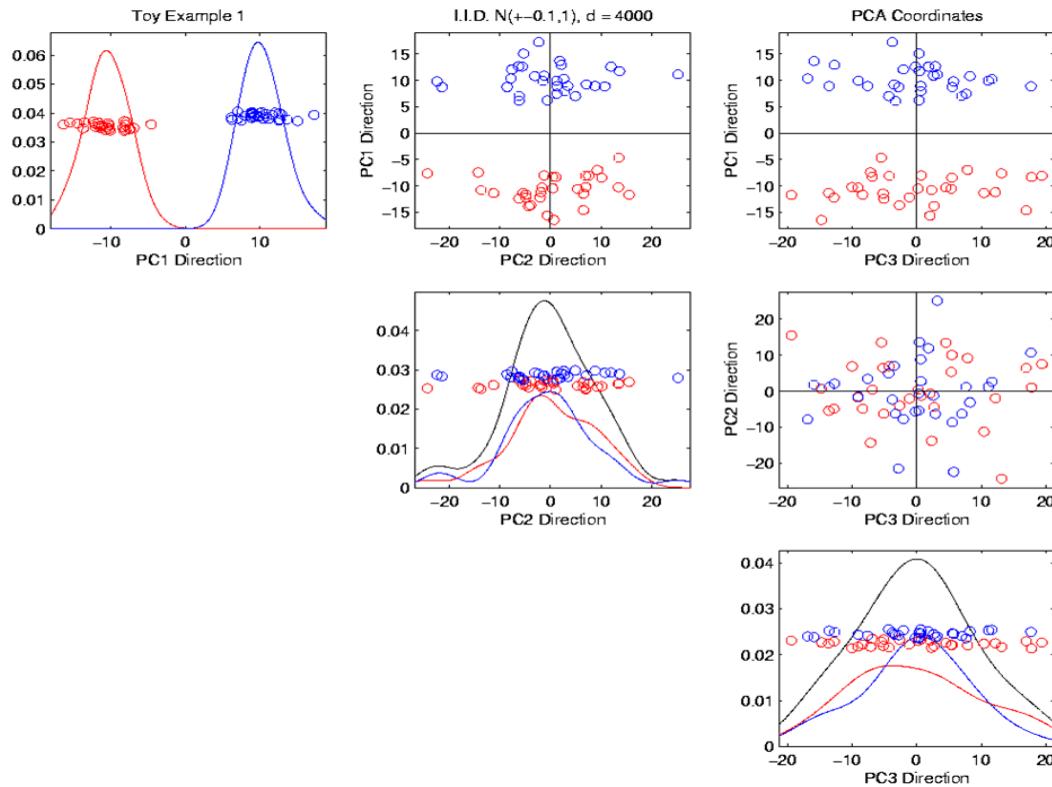
- Now with much higher dimension, # genes = 4000
- 2 clusters with slight mean difference in each variable $N(\pm 0.1, 1)$

Gene by Gene view



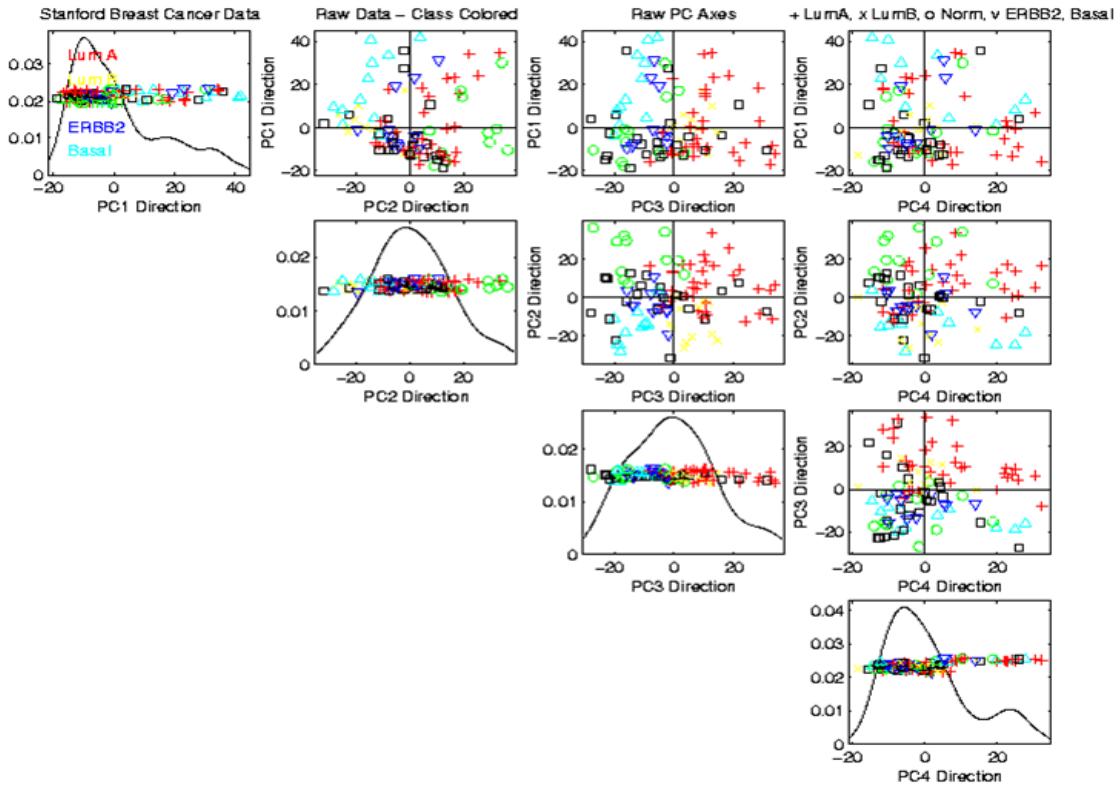
2 clusters nearly the same

PCA view

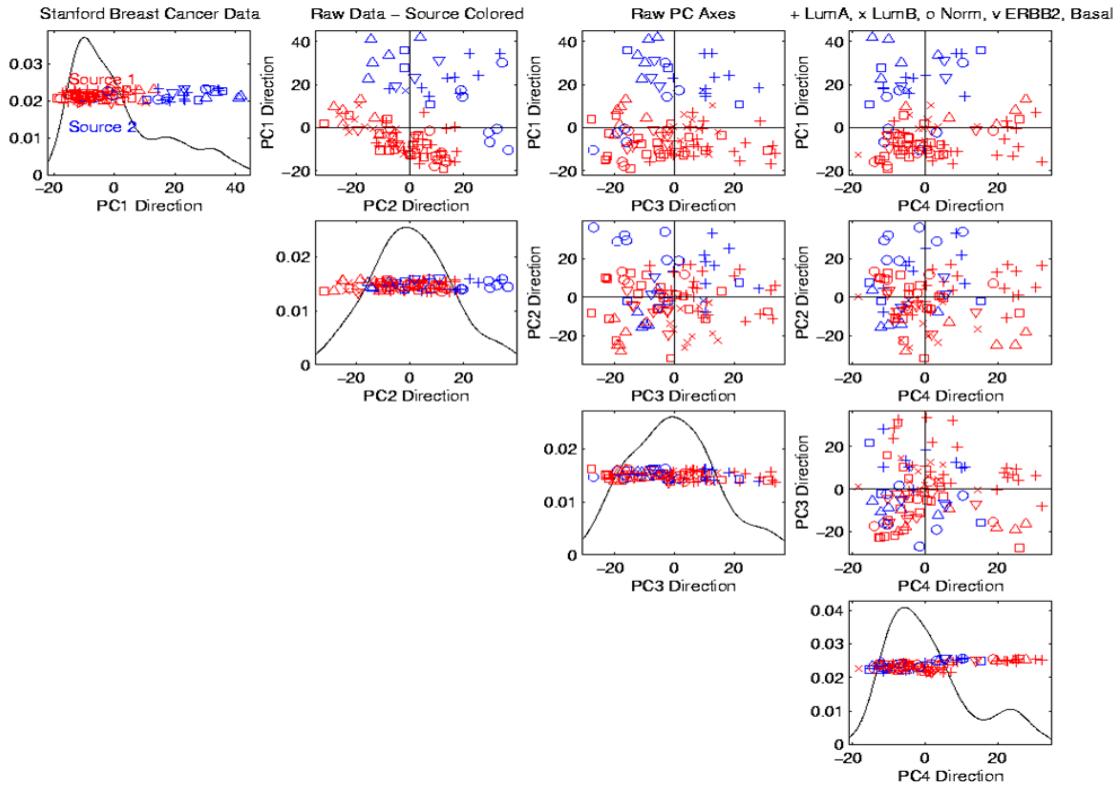


Huge difference in the 1st PC direction
Alternative views can show more

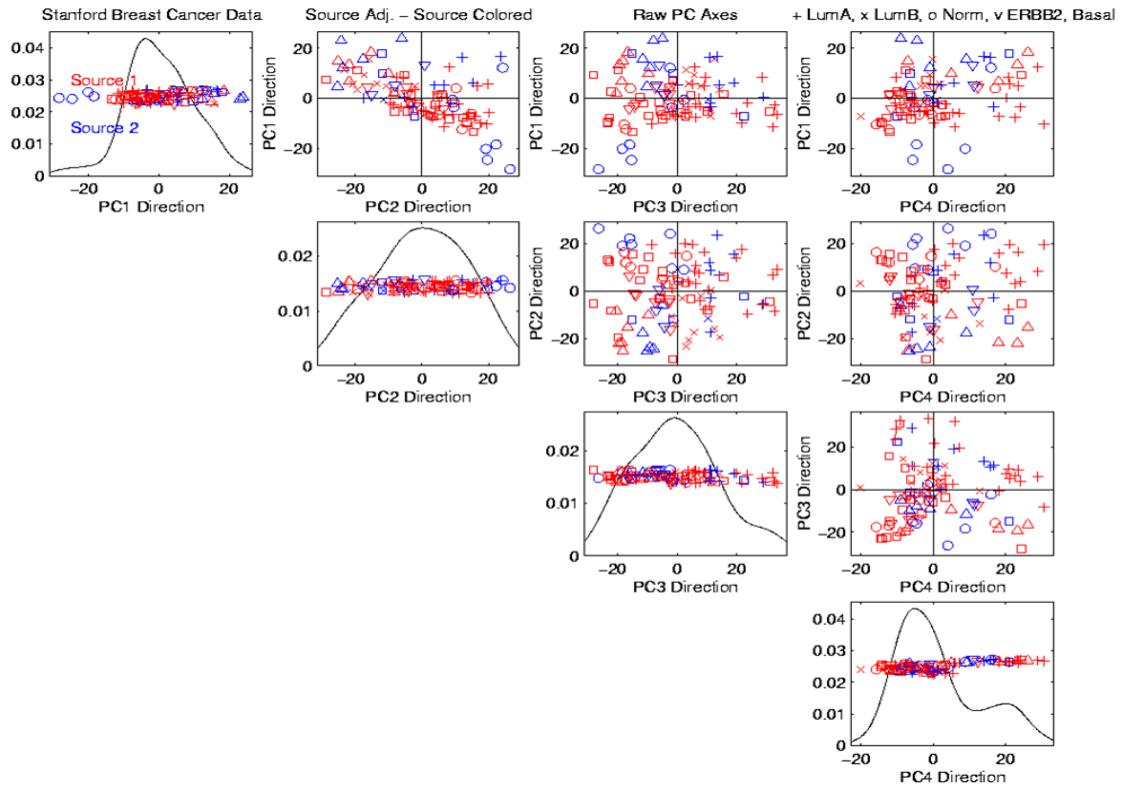
Breast Cancer Data



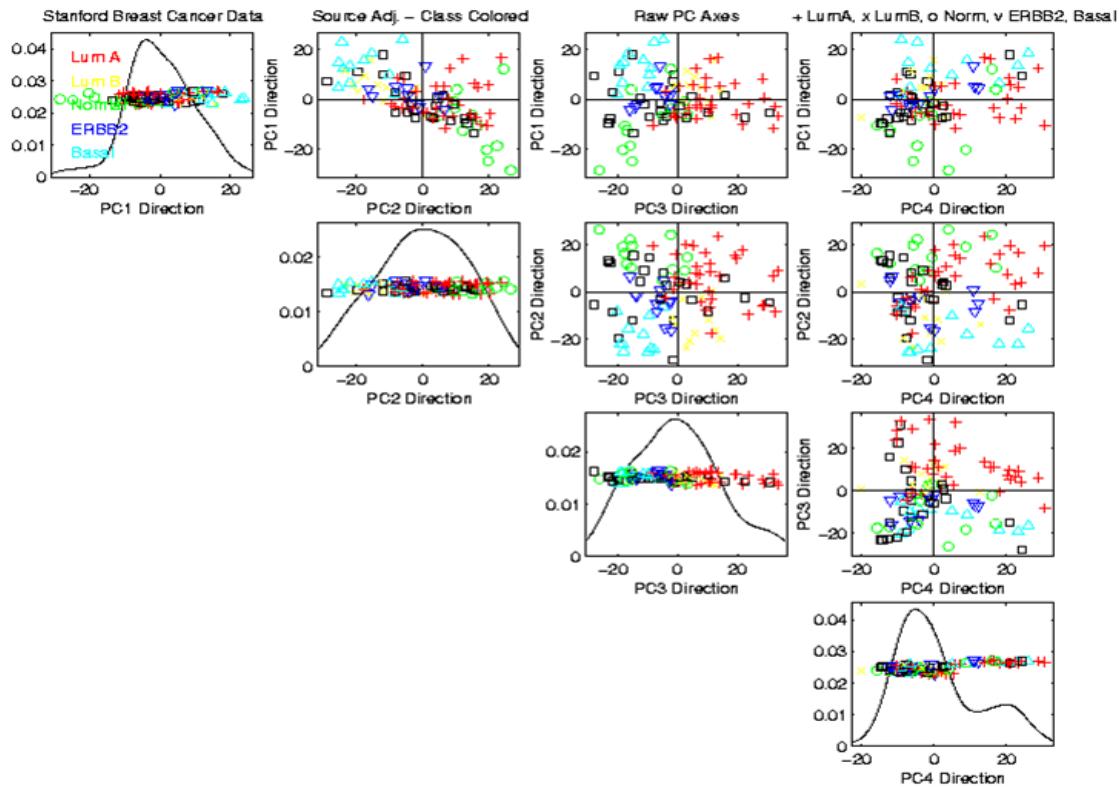
Breast Cancer Data



Breast Cancer Data



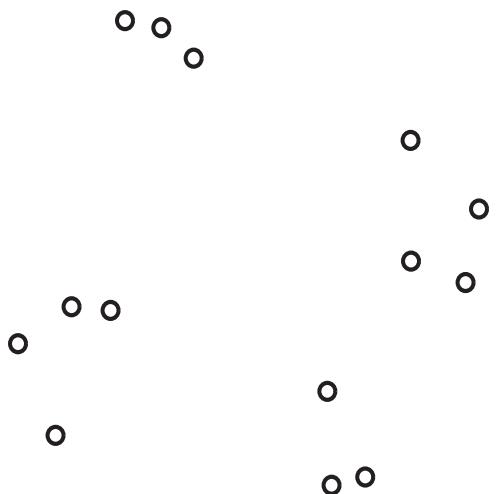
Breast Cancer Data



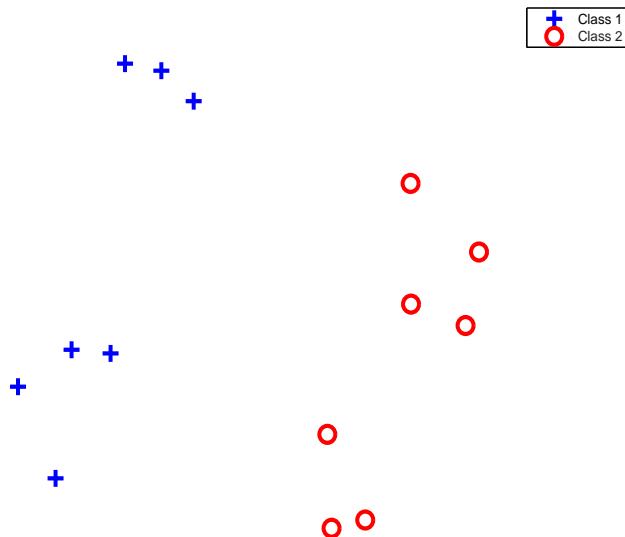
Classification

- Class 1: $\mathbf{X}_{d \times m}^1$, Class 2: $\mathbf{X}_{d \times n}^2$
- Example: microarray from cancer / non-cancer patients.
- Use training data from Class 1 & Class 2.
- Develop a rule for assigning new data to a class label.
- Based on measurements, want to predict if new patients are Healthy or Ill

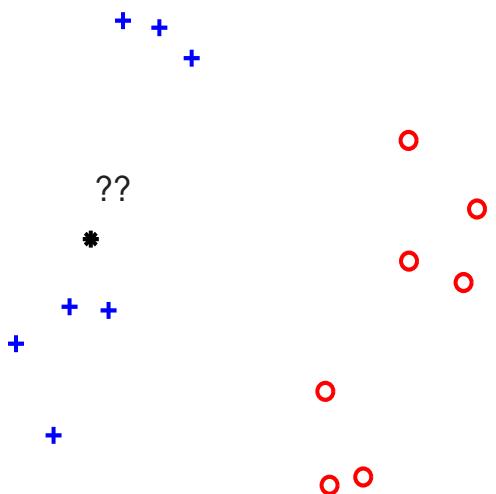
Example



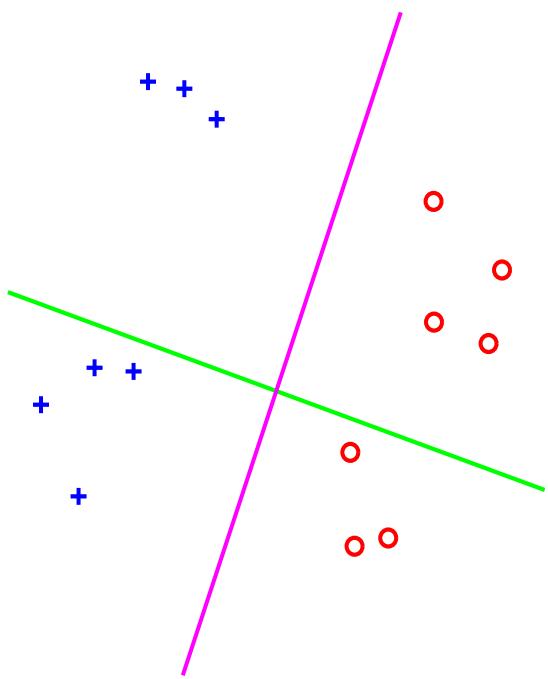
Example



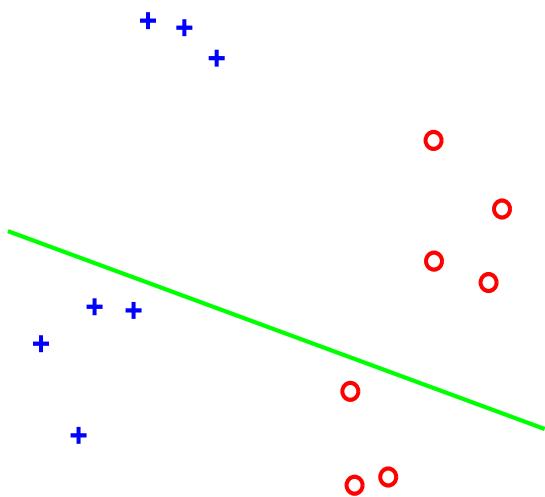
Example



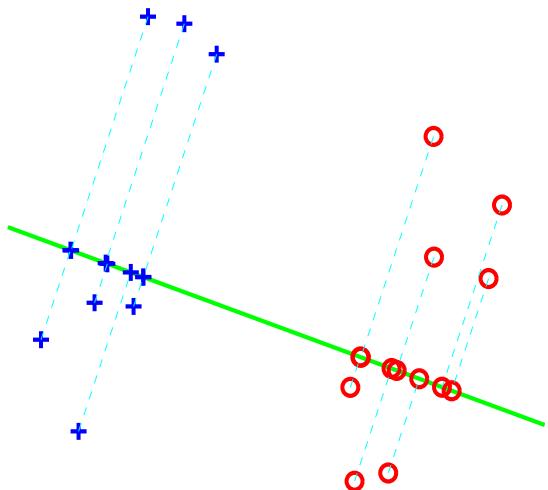
Example



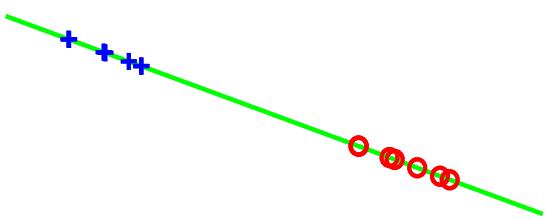
Example



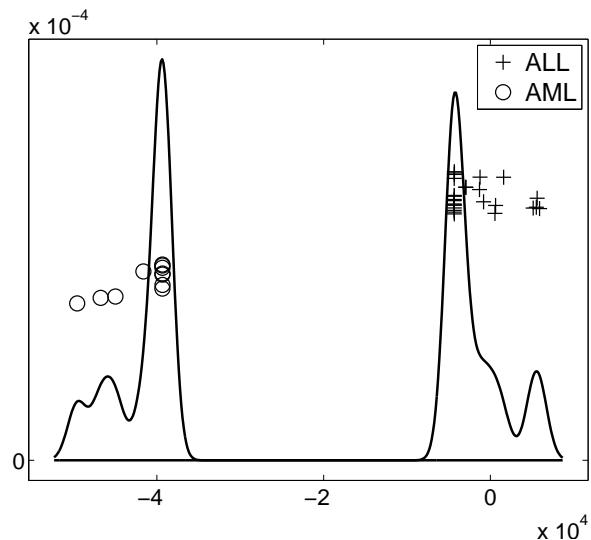
Example



Example



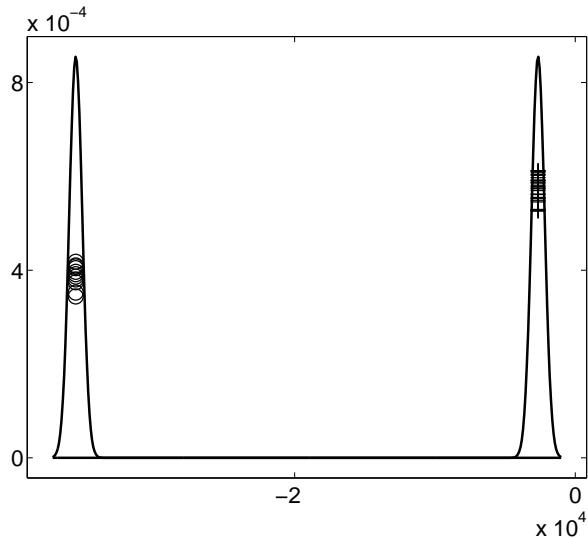
Classification visualization



$n = 38$ and $p = 7129$

Projections of leukemia microarray data onto the SVM direction vector.

Classification visualization



$n = 38$ and $p = 7129$

Projections of leukemia microarray data onto the MDP direction vector (defined only for HDLSS data).

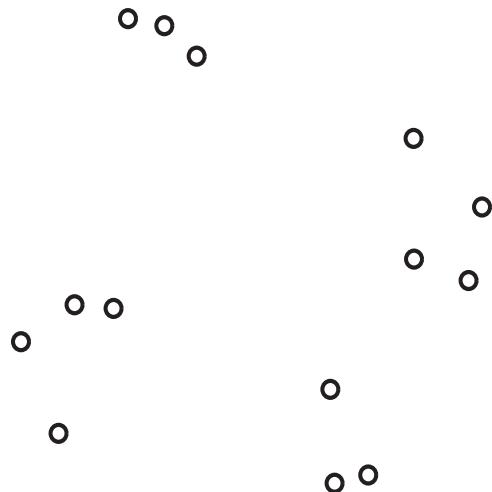
Cluster Analysis

Goal: we want to see if some natural groups (classes, clusters) of individuals exists.

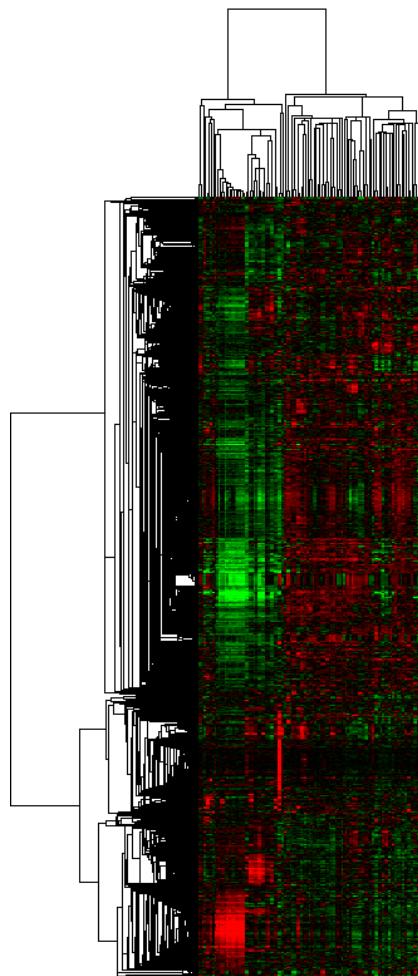
The groups should be as homogeneous as possible and the differences between the groups as large as possible.

- Choice of proximity measure
- Choice of grouping algorithm

Cluster Analysis Example



Example: Ovarian Cancer Data



$n = 108$ and $p = 28,475$

Goal

- Understand why and how multivariate methods are used
- Understand when a multivariate method should be used
- Understand some of the theory behind the methods
- Gain experience applying multivariate methods and interpreting results